

## ABSTRACT

Title of dissertation: PLANETARY SCALE DATA STORAGE

Benjamin Bengfort  
Doctor of Philosophy, 2018

Dissertation directed by: Professor Peter J. Keleher  
Department of Computer Science

The success of virtualization and container-based application deployment has fundamentally changed computing infrastructure from dedicated hardware provisioning to on-demand, shared clouds of computational resources. One of the most interesting effects of this shift is the opportunity to localize applications in multiple geographies and support mobile users around the globe. With relatively few steps, an application and its data systems can be deployed and scaled across continents and oceans, leveraging the existing data centers of much larger cloud providers.

The novelty and ease of a global computing context means that we are closer to the advent of an Oceanstore, an Internet-like revolution in personalized, persistent data that securely travels with its users. At a global scale, however, data systems suffer from physical limitations that significantly impact its consistency and performance. Even with modern telecommunications technology, the latency in communication from Brazil to Japan results in noticeable synchronization delays that violate user expectations. Moreover, the required scale of such systems means that failure is routine.

To address these issues, we explore consistency in the implementation of distributed logs, key/value databases and file systems that are replicated across wide areas. At the core of our system is hierarchical consensus, a geographically-distributed consensus algorithm that provides strong consistency, fault tolerance, durability, and adaptability to varying user access patterns. Using hierarchical consensus as a backbone, we further extend our system from data centers to edge regions using federated consistency, an adaptive consistency model that gives satellite replicas high availability at a stronger global consistency than existing weak consistency models.

In a deployment of 135 replicas in 15 geographic regions across 5 continents, we show that our implementation provides high throughput, strong consistency, and resiliency in the face of failure. From our experimental validation, we conclude that planetary-scale data storage systems can be implemented algorithmically without sacrificing consistency or performance.

# PLANETARY SCALE DATA STORAGE

by

Benjamin Bengfort

Dissertation submitted to the Faculty of the Graduate School of the  
University of Maryland, College Park in partial fulfillment  
of the requirements for the degree of  
Doctor of Philosophy  
2018

Advisory Committee:  
Professor Peter J. Keleher, Chair/Advisor  
Professor Dave Levin  
Professor Amol Deshpande  
Professor Daniel Abadi  
Professor Derek C. Richardson

© Copyright by  
Benjamin Bengfort  
2018



## Preface

If needed.

## Foreword

If needed.

## Dedication

To Irena and Henry.



## Acknowledgments

I could not have done this alone.

## Table of Contents

Preface	ii
Foreword	iii
Dedication	iv
Acknowledgements	v
List of Tables	viii
List of Figures	ix
List of Abbreviations	x
1 Introduction	1
2 Challenges and Motivations	6
2.1 A New Application Development Paradigm . . . . .	8
2.2 Building Geo-Replicated Services . . . . .	11
2.3 Requirements for Data Systems . . . . .	12
2.4 System Architecture . . . . .	12
2.5 Conclusion . . . . .	13
3 Hierarchical Consensus	14
3.1 Overview . . . . .	16
3.2 Consensus . . . . .	19
3.2.1 Root Consensus . . . . .	22
3.2.2 Delegation . . . . .	24
3.2.3 Epoch Transitions . . . . .	26
3.2.4 Fuzzy Handshakes . . . . .	29
3.2.5 Subquorum Consensus . . . . .	32
3.2.6 Client Operations . . . . .	33
3.3 Consistency . . . . .	34

3.3.1	Globally Consistent Logs . . . . .	34
3.4	Fault Tolerance . . . . .	40
3.4.1	Failures . . . . .	42
3.4.2	Obligations Timeout . . . . .	43
3.4.3	The Nuclear Option . . . . .	43
3.5	Performance Evaluation . . . . .	45
3.6	Conclusion . . . . .	53
4	Federated Consistency . . . . .	55
4.1	Overview . . . . .	55
4.2	Eventual Consistency . . . . .	55
4.2.1	Consistency Failures . . . . .	55
4.3	Integration . . . . .	56
4.3.1	Communication Integration . . . . .	56
4.3.2	Consistency Integration . . . . .	56
4.4	Performance Evaluation . . . . .	56
5	System Implementation . . . . .	57
5.1	System Model . . . . .	58
5.2	Applications . . . . .	59
5.2.1	Distributed Log . . . . .	59
5.2.2	Key-Value Database . . . . .	59
5.2.3	File System . . . . .	59
5.3	Consistency Model . . . . .	59
5.4	Raft . . . . .	60
5.5	Conclusion . . . . .	63
6	Adaptive Consistency . . . . .	65
6.1	Bandit-Based Approaches . . . . .	65
6.1.1	Rewards Function . . . . .	65
6.2	Access Temperature Approaches . . . . .	65
7	Related Work . . . . .	66
7.1	Hierarchical Consensus . . . . .	66
7	Conclusion . . . . .	69
A	Formal Specification . . . . .	70
	Bibliography . . . . .	71

## List of Tables

3.1	HC Failure Categories . . . . .	41
5.1	Parameterized Timeouts of Raft Implementation . . . . .	62

## List of Figures

1.1	Global Data Centers of Cloud Providers . . . . .	2
2.1	Distributed Architectures . . . . .	11
2.2	Global Architecture . . . . .	13
3.1	A 12x3 Hierarchical Consensus Network Topology . . . . .	18
3.2	HC Operational Summary . . . . .	20
3.3	Delegated Votes . . . . .	23
3.4	Ordering of Epochs and Terms in Root and Subquorums . . . . .	28
3.5	Epoch Transition: Fuzzy Handshakes . . . . .	30
3.6	Grid Consistency: A Sequential Log Ordering . . . . .	36
3.7	Sequential Event Ordering in HC . . . . .	38
3.8	Event Ordering with Remote Writes in HC . . . . .	39
3.9	Scaling Consensus HC vs. Raft . . . . .	47
3.10	HC Throughput vs. Workload in the Wide Area . . . . .	48
3.11	HC Cumulative Latency Distribution . . . . .	49
3.12	Sawtooth Graph . . . . .	51
3.13	HC Fault Repartitioning . . . . .	52

## List of Abbreviations

EC	Eventual consistency
HC	Hierarchical consensus

## Chapter 1: Introduction

Eighteen years ago the Oceanstore paper [1] presented a vision for a data utility infrastructure that spanned the globe. The economic model was that of a cooperative utility provided by a confederation of companies that could buy and sell capacity to directly support their users, with regional providers like airports and cafes installing servers to enhance performance for a small dividend of the utility. This economic model meant that Oceanstore’s requirements centered around an untrusted infrastructure to support nomadic data: connectivity, security, durability, and location agnostic storage. To meet these requirements, the Oceanstore architecture was composed of two tiers: pools of byzantine quorums that made localized consistency and placement decisions, along with an optimistic dissemination tree layer that moved data between quorums as correctly as possible without providing guarantees. This architecture, along with a reliance on encryption and key-based access control, could facilitate grid computing storage, a truly decentralized and independent participation of heterogenous computational resources across the globe [2].

Unforeseen by Oceanstore, however, was a fundamental shift in how companies and users accessed computing infrastructure. Improvements in virtualization

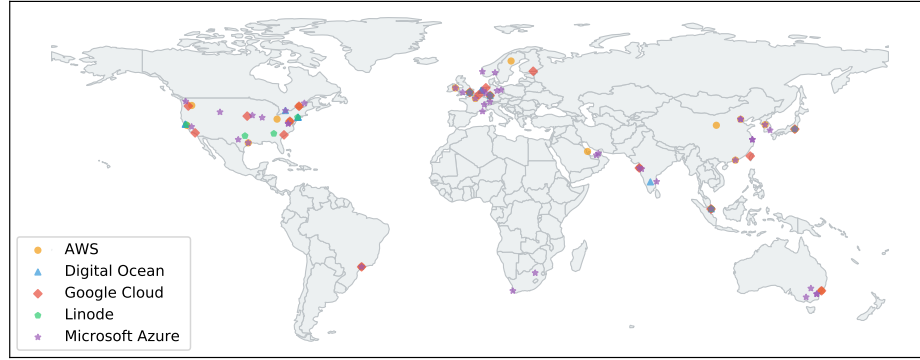


Figure 1.1: Data center locations of popular cloud providers span the globe.

management [3] and later container computing [4] allowed big internet companies to easily lease their unused computational resources and disk capacity to application developers, making cloud computing [5] rather than independent hardware purchasing and hosting the norm. Furthermore, from smarter phones to tablets and netbooks that do not have large disk capacities, user devices have become increasingly mobile; from photos to email and contacts, most user data is now stored in cloud silos. The cloud economy means that there exists a *trusted infrastructure* of virtual resources that span globe, provisioned by a single provider as shown in Figure 1.1.

So what does this mean for the requirements and assumptions of Oceanstore? First, the strict requirements for security that meant per-user encryption and byzantine agreement between untrusted servers can be relaxed to application and transport-level encryption and non-byzantine consensus supported by authenticated communication. Second, the requirements for performance have dramatically increased as ubiquitous computing has become the norm and as more non-human users are par-



ticipating in networks. Increased capacity, however, cannot come at the cost of correctness or consistency, and the increased rate of requests means that asynchronous commits and conflict resolution become far more difficult. For these reasons, we believe that a vision for an Oceanstore today would focus on *consistency* rather than security.

Consistent behavior in distributed systems becomes increasingly complex to implement and reason about as the system size grows and requires increased coordination. In 2021, Cisco forecasts over 25 billion devices will contribute to 105,800 GBps of global internet traffic, 26% of which will be file sharing and application data, and 51% of which will originate from machine to machine-to-machine applications [6]. New types of networks including sensor networks, smart grid solutions, self-driving vehicle networks, and an internet of things will mean an update model with many publishers, few subscribers, and increasingly distributed accesses. To support this growth and facilitate speed, traffic is consistently moving closer to the edge; Cisco predicts that cross-country delivery will drop from 58% of traffic in 2016 to 41% in 2021 and that metro delivery will grow from 22% to 35%. Localization means that the cloud will be surrounded by a fog of devices that participate in systems by contributing data storage and computation to an extent greater than access-oriented clients might.

Based on these trends and inspired by the work of Oceanstore, we propose that a consistent planetary-scale data storage system would be made up of a two tier architecture of both cloud and fog infrastructure. The first tier, in the cloud, would be a strong consistency, fault tolerant and highly resilient geo-replicated con-

sensus backbone: *hierarchical consensus*. The second tier, via the fog, would be a high availability heterogeneous network with a hybrid consistency model: *federated consistency*. Such a system would be difficult to manually manage, therefore the system would also have to automatically monitor and adapt to changes in access patterns and node and network availability during runtime. We believe that the combination of hierarchical consensus, federated consistency, and adaptive monitoring lay out a foundation for truly large scale data storage systems that span the planet.

The contributions of this dissertation are therefore as follows:

1. We present the design, implementation, and evaluation of hierarchical consensus, a consensus protocol that can scale to dozens or hundreds of replicas across the wide area.
2. We also investigate the design and implementation of federated consistency, a hybrid consistency model that allows strong, consensus-based systems to integrate with eventually-consistent, highly available replicas and evaluate it in a simulated heterogeneous network.
3. We show the possibilities for machine learning-based system adaptation with a reinforcement learning approach to anti-entropy synchronizations based on accesses.
4. We validate our system by describing the implementation of a planetary-scale key-value data store and file system using both hierarchical consensus and federated consistency.

The rest of this dissertation is organized as follows. In the next chapter we will more thoroughly describe the motivations and challenges of building geo-replicated data systems as well as explore case-studies of existing systems. Next, we will focus on the core backbone of our system: hierarchical consensus and describe a globally fault tolerant approach to managing accesses to objects in the wide area. Using hierarchical consensus as a building block, we will next describe federated consistency and how a hybrid, heterogenous consistency model in the fog interacts with the cloud consensus tier. At this point we will have enough background to introduce our system implementation and describe our file system and key-value store. From there, we will explore learning systems that monitor and adapt the performance of the system at runtime, before concluding with related work and a discussion of our future research.

## Chapter 2: Challenges and Motivations

Many of the world's most influential companies grew from the ashes of the dot-com bubble of the 1990s, which paid for an infrastructure of fiber-optic cables, giant server farms, and research into mobile wireless networks [7]. As these companies filled market voids in eCommerce, search, and social networking, they created new database technologies to leverage the potential of underused computational resources and low latency/high bandwidth networks that connected them, eschewing more mature systems that were developed with resource scarcity in mind [8, 9]. What followed was the rise and fall of NoSQL data systems, a microcosm of the proceeding era of database research and development [10]. Although there are a lot of facets to the story of NoSQL, what concerns us most is the use of NoSQL to create geographically distributed systems, as these systems paved the way to the large-scale storage systems in use today.

The first phase of distributed NoSQL systems was the creation of highly available, sharded systems intended to meet the demand of increasing numbers of clients. Commercially, these types of systems include Dynamo [11] and BigTable [12], which in turn spawned open source and academic derivatives such as Cassandra [13] and HBase [14]. Although these systems did support large number of accesses, they

achieved their availability by relaxing consistency, which many applications found to be intolerable. The second phase was, therefore, a return to stronger consistency, even at the cost of decreased performance or expensive engineering solutions. Again, commercial systems led the way with Megastore [15] and Spanner [16] along with academic solutions such as MDCC [17] and Calvin [18, 19]. Part of this realignment was a reconsidering of the base assumption that drove the NoSQL movement expressed in the CAP theorem [20], primarily that the lines between availability, partition tolerance, and consistency may not be as strictly drawn as previously theorized [21, 22]. This has led to the beginning of a third phase, the return of SQL, as the lessons learned during the glut of computational resources are applied to more traditional systems. As before, both commercial systems, such as Aurora [23] and Azure SQL [24], and open source systems such as Vitess [25] and CockroachDB [26] are playing an important role in framing the conversation about consistency in this phase.

This brief and limited description of the history of NoSQL distributed systems serves to set the tone for two primary points. First, designing distributed systems to support a large number of users across large geographies entails a number of challenges and trade-offs due to physical limitations that no single architecture has been able to fully cover. Second, there exists commercial and practical motivations to build such systems and these motivations have been the impulse toward academic research. With this backdrop in mind, we explore in this chapter the question of whether a planetary-scale data system is really necessary, and the challenges that we face when designing such systems.

## 2.1 A New Application Development Paradigm

The launch of the augmented reality game Pokémon GO in the United States was an unmitigated disaster [27]. Due to extremely overloaded servers from the release’s extreme popularity, users could not download the game, login, create avatars, or find augmented reality artifacts in their locales. The company behind the platform, Niantic, scrambled quickly, diverting engineering resources away from their feature roadmap toward improving infrastructure reliability. The game world was hosted by a suite of Google Cloud services, primarily backed by the Cloud Datastore [28], a geographically distributed NoSQL database. Scaling the application to millions of users therefore involved provisioning extra capacity to the database by increasing the number of shards as well as improving load balancing and autoscaling of application logic run in Kubernetes [29] containers.

Niantic’s quick recovery is often hailed as a success story for cloud services and has provided a model for elastic, on demand expansion of computational resources. A deeper examination, however, shows that Google’s global high speed network was at the heart of ensuring that service stayed stable as it expanded [30]. The original launch of the game was in 5 countries – Australia, New Zealand, the United States, the United Kingdom, and Germany; however the success of the game meant worldwide demand, and it was subsequently expanded to over 200 countries starting with Japan [31]. Unlike previous games that were restricted with region locks [32], Pokémon GO was a truly international phenomenon and Niantic was determined to allow international interactions in the game’s feature set, interaction which relies on

Google’s unified international architecture and globally distributed databases.

In the brief history of NoSQL that started this chapter, we showed that the growth of database systems distributed across the wide-area started with large internet companies like Yahoo, Google, and Amazon but quickly led to academic investigations. One reason that the commercial systems enjoyed this academic attention was that at the time, the unique scale of their usage proved the motivation behind their architecture, however their success has meant that these types of scales are no longer limited to huge software systems. Instead, stories such Niantic’s deployment are becoming common and medium to large applications now require developers to increasingly reason about how data is distributed in the wide area, different political regions, and replicated for use around the world.

Consider the following companies and applications from large to small that have international audiences. Dropbox has users in over 180 countries and is supported in 20 languages, maintaining offices in 12 locations from Herzliya to Sydney [33]. Slack serves 9 million weekly active users around the world and has 8 offices around the world, prioritizing North America, Europe, and Pacific regions [34]. WeWork provides co-working space in 250+ international applications and uses an app to manage global access and membership [35]. Tile has sold 15 million of its RFID trackers worldwide and locates 3 million unique items a day in 230 countries [36]. Trello, a project management tool, has been translated into 20 languages and has 250 million world-wide users in every country except Tuvalu, their international roll-out focused on marketing and localization [37]. Runkeeper [38] and DarkSky [39] are iOS and Android apps that have millions of global users and struggled to make their

services available in other countries, but benefitted from international app stores. Signal and Telegraph, encrypted messaging apps have grown primarily in countries at the top of Transparency International’s Corruption Perception Index [40].

None of the applications described above necessarily have geography-based requirements in the same way that an augmented reality or airline reservations application might have, just a large number of users who regularly use the app from a variety of geographic locations. Web developers are increasingly discussing and using container based approaches both for development and small-scale production, web frameworks have built in localization tools that are employed by default, and services are deployed on autoscaling cloud platforms from the start. The new application development paradigm, even for small applications, is to build with the thought that your application will soon be scaling across the globe.

To address this paradigm, cloud service providers have expanded their offerings to include usage of their distributed data stores to application developers. The problem is that distributed systems like Dynamo, BigTable, Megastore, and Spanner will all be designed for in-house services with a data model that supported replication and eventual consistency, but makes it difficult for developers to reason about behavior, as we will see in the next section. Not only is there a need for strong consistency semantics, data-location awareness, and geo-replication in distributed data storage systems, there is also the need for a familiar and standardized storage API.



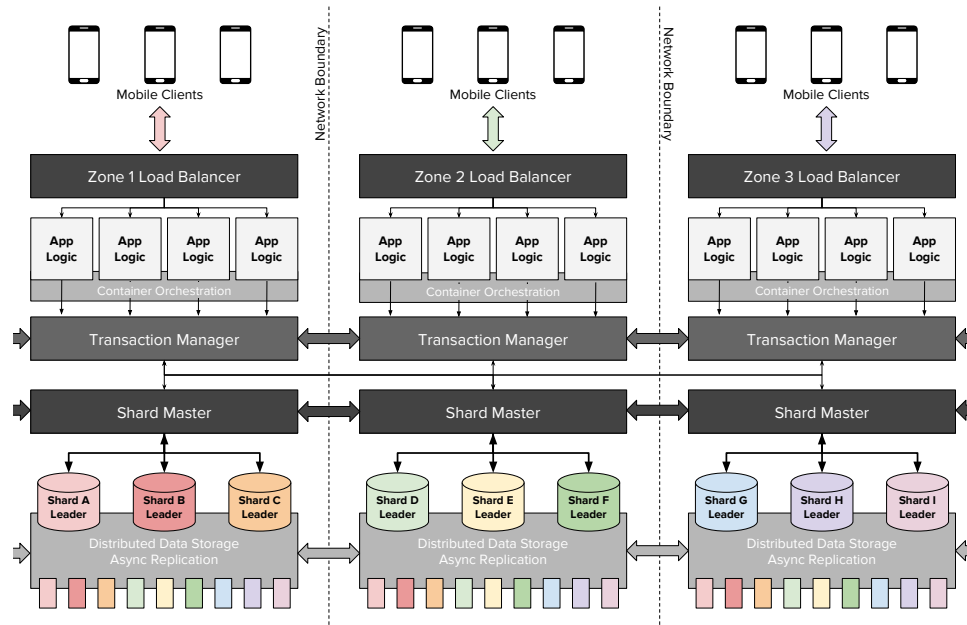


Figure 2.1: A standard architecture of a distributed application.

## 2.2 Building Geo-Replicated Services

Given non-geographic requirements for building a distributed system, a common architecture that provides high write throughput with consistent replication

During the first phase of

Paxos as the basis for high performance data store [41]

What they didn't do and what we do better.

BigTable & Spanner

S3

Aurora, Cosmos, & Cockroach DB

Approaches: sharding, independent objects, buckets, slow reads

## 2.3 Requirements for Data Systems

Failure is common Disk failure/replica failure (ODS) Network failure (when repaired, replica comes back online) Unreliability: messages have highly variable latency, out of order messaging Partitions, part of the system cannot speak to the rest of the system

In geo-systems large latency is not the issue! There is a physical limit to message traffic Writes must be applied in order, reads can reason about staleness Access patterns are also location-dependent

Durability Normally 3 disk replication ensures 2 failures We need to ensure zone+1 disk failures (re Aurora) User-specific data should be accessible everywhere

Fault-Tolerance & availability System should still be available if nodes fail Should be available if zones fail (e.g. hurricane or disaster) at higher latency System should be available at lower consistency if even one replica is available

Adaptability Should respond to changes in user access patterns Should be able to add and remove nodes from the system Should scale with more regions and more replicas

## 2.4 System Architecture

Describe the architecture of tier 1: HC, tier 2: Federated Fog

Describe the consistency guarantees we claim

Base application is a key/value store

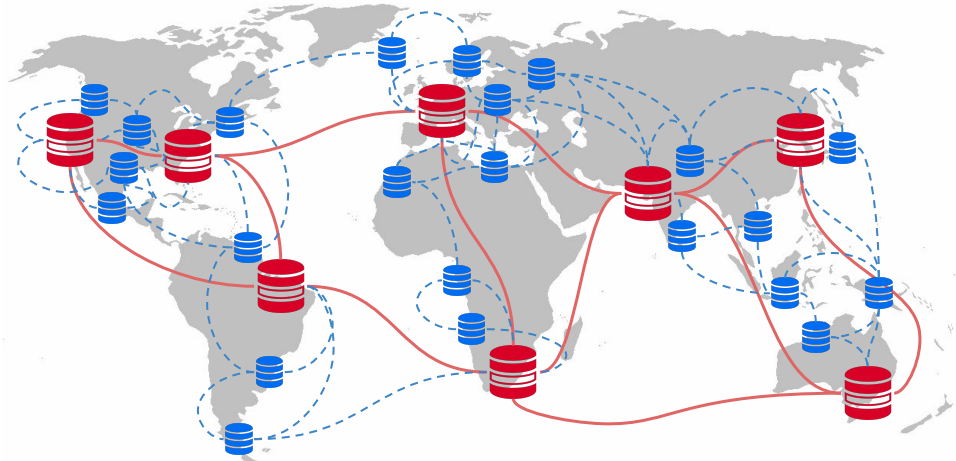


Figure 2.2: A global architecture composed of a core backbone of hierarchical consensus replication (red) and a fog of heterogeneous, federated consistency replicas (blue).

File-system is built upon the key/value store as an object FS

In Figure [2.2](#)

## 2.5 Conclusion

1. Global-scale applications are becoming the new norm
2. Today's globally distributed data systems are high level and require developers to deeply consider consistency and localization semantics.
3. The challenges are around the physical limitations of geographic networks: the speed of light in this case; e.g. latency and outages.

## Chapter 3: Hierarchical Consensus

The backbone of our planetary scale data system is *hierarchical consensus*. Hierarchical consensus provides a strong consistency foundation that totally orders critical accesses and arbitrates the eventual consistency layer in the fog, which raises the overall consistency of the system. To be effective, an externalizable view of consistency ordering must be available to the entire system. This means that strong consistency must be provided across geographic links rather than provided as localized, independent decision making with periodic synchronization. The problem that hierarchical consensus is therefore designed to solve is that of distributed consensus.

Solutions to distributed consensus primarily focus on providing high throughput, low latency, fault tolerance, and durability. Current approaches [18, 42–46] usually assume a small number of replicas, each centrally located on a highly available, powerful, and reliable host. These assumptions are justified by the environments in which they run: highly curated environments of individual data centers connected by dedicated networks. Although replicas in these environments may participate in global consensus, our data model requires us to accommodate replicas with heterogeneous capabilities and usage modalities. Widely distributed replicas might have neither high bandwidth nor low latency and might suffer partitions of varying du-

rations. Such systems of replicas might also be dynamic in membership, in relative locations, and have relative workloads. Most importantly, to provide a backbone for a planetary scale data system, the consistency backbone must scale to include potentially hundreds of replicas around the world.

As a result, straightforward approaches of running variants of Paxos [47], ePaxos [42], or Raft [48] across the wide area, even for individual objects will perform poorly for several reasons. First, distance (in network connectivity) between consensus replicas and the most active replicas decrease the performance of the entire system, consensus is only as fast as the final vote required to make a decision, even when making thrifty requests. Second, network partitions are common, which cause consensus algorithms to fail-stop [49] if they cannot receive a majority, a criticism that is often used to justify eventual consistency systems for high availability. Finally, the fault tolerance of small quorum algorithms can be disrupted by a small number of unreliable hosts and given the scale of the system and the heterogenous nature of replicas, the likelihood of individual failure is so high so as to be considered inevitable.

We propose another approach to building large systems. Rather than relying on a few replicas to provide consensus to many clients, we propose to run a consensus protocol across replicas running at or near all of these locations. The key insight is that large problem spaces can often be partitioned into mostly disjoint sets of activity without violating consistency. We exploit this decomposition property by making our consensus protocol hierarchical and individual consensus groups fast by ensuring they are small. We exploit locality by building subquorums from colocated

replicas, and locating subquorums near clients they serve.

In this chapter we describe hierarchical consensus, a two-tiered consensus structure that allows high throughput, localization, agility, and linearizable access to a shared namespace. We show how to use *delegation* to build large consensus groups that retain their fault tolerance properties while performing like small groups. We describe the use of *fuzzy epoch transitions* to allow global re-configurations across multiple consensus groups without forcing them into lockstep. Finally, we describe how we reason about consistency by describing the structure of grid consistency.

### 3.1 Overview

Hierarchical Consensus (HC) is an implementation and extension of Vertical Paxos [50–52] designed to scale to hundreds of nodes geo-replicated around the world. Vertical Paxos divides consensus decisions both horizontally, as sequences of consensus instances, and vertically as individual consensus decisions are made. Spanner [16], MDCC [17], and Calvin [18], can all be thought of as implementations of Vertical Paxos, in that they shard the namespace of the objects they manage into individual consensus groups (the vertical division). In these cases, however, sharding does not allow for inter-object dependence (the horizontal division) without either a management quorum which is either not geo-replicated, suffers from the same problems in scaling, or without the use of extremely accurate timestamps. The challenge is therefore in building a multi-group coordination protocol that configures and mediates the subquorums with the same level of consistency and fault tolerance

of the entire system.

Hierarchical consensus therefore organizes *all* participating replicas to participate in a root quorum as shown in Figure 3.1. The root quorum guarantees correctness by pivoting the overall system through two primary functions. First, the root quorum reconfigures subquorum memberships on replica failures and system membership changes (allocating hot-spares as needed). Second, the root quorum adjusts the mapping of the object namespace to the underlying partitions. Much of the system's complexity comes from handshaking between the root quorum and the lower-level subquorums during reconfigurations.

These handshakes are made easier, and much more efficient, by using *fuzzy transitions*. Fuzzy transitions allow individual subquorum to move through reconfiguration at their own pace, allowing portions of the system to transition to decisions made by the root quorum before others. Given out heterogenous, wide-area environment, forcing the entire system to transition to new configurations in lockstep would be unacceptably slow. Fuzzy transitions also ensure that there is no dedicated shard-master that has to synchronize all namespace allocations: at the cost of possibly multiple redirections, clients can be redirected by any member of the root quorum to replicas who should be participating in consensus decisions for the requested objects.

Fuzzy transitions ensure that root quorum decisions need not be timely since those decisions do not disrupt accesses of clients. Though root quorum decisions are rare with respect to the throughput of accesses inside the entire system, they still do require the participation of all members of the system, which could lead

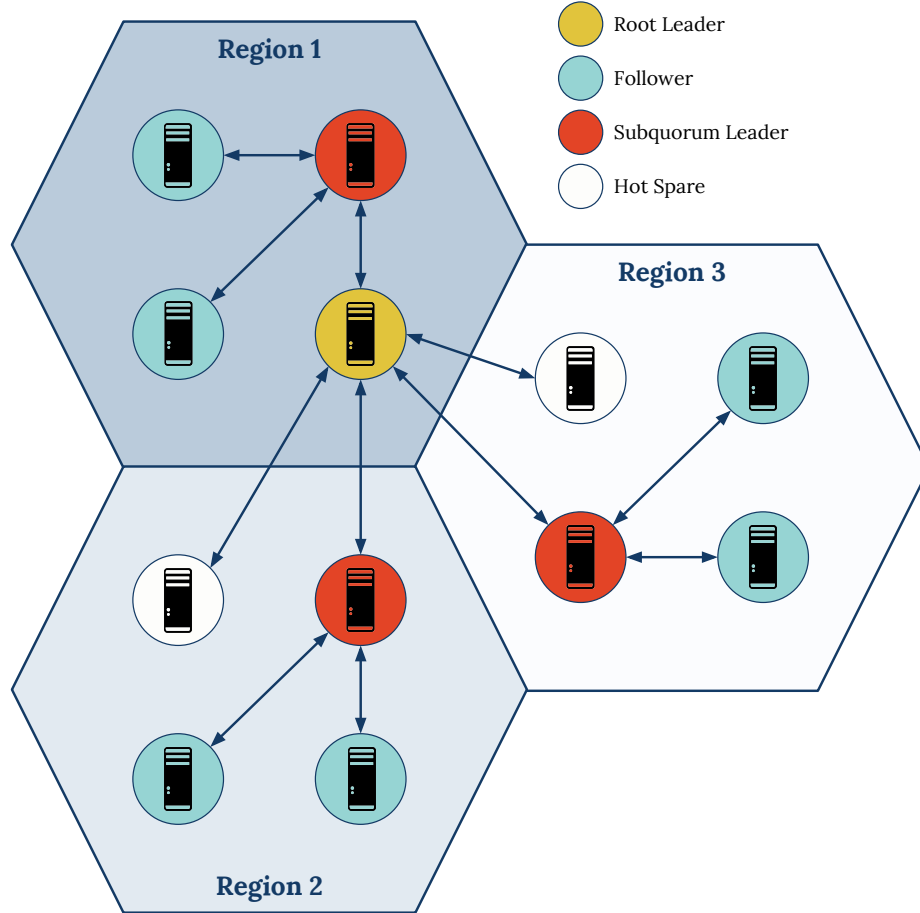


Figure 3.1: A simple example of an HC network composed of 12 replicas with size 3 subquorums. Each region hosts its own subquorum and subquorum leader, while the subquorum leaders delegate their votes to the root quorum, whose leader is found in region 1. This system also has 2 hot spares that can be used to quickly reconfigure subquorums that experience failures. The hot spares can either delegate their vote, or participate directly in the root quorum.



to extremely large quorum sizes, and therefore extremely slow consensus operations that may be extremely sensitive to partitions. Because all subquorums make disjoint decisions and because all members of the system are part of the root quorum, we propose a safe relaxation of the participation requirements for the root quorum such that subquorum followers can *delegate* their root quorum votes to their leader. Delegation ensures that only a small number of replicas participate in most root quorum decisions, though decisions are made for the entire system.

In brief, the resulting system is local, in that replicas serving clients can be located near them. The system is fast because individual operations are served by small groups of replicas, regardless of the total size of the system. The system is nimble in that it can dynamically reconfigure the number, membership, and responsibilities of subquorums in response to failures, phase changes in the driving applications, or mobility among the member replicas. Finally, the system is consistent, supporting the strongest form of per-object consistency without relying on special-purpose hardware [16, 53–56].

A complete summary of hierarchical consensus is described in Figure 3.2.

## 3.2 Consensus

The canonical distributed consensus used by systems today is Paxos [47, 57]. Paxos is provably safe and designed to make progress even when a portion of the system fails. Raft [48] was designed not to improve performance, but to increase understanding of consensus behavior to better allow efficient implementations. HC

Root Management	Delegated Votes	“Nuclear” Option
	<p>Discussed in §2.1</p> <p><b>Root Leader</b></p> <ul style="list-style-type: none"><li>• Broadcast command to all replicas.</li><li>• Resolves conflicts (<b>q,t</b>) by selecting the delegation with highest term.</li><li>• If current vote count is a majority, begin epoch transition.</li></ul> <p><b>Root Delegates</b></p> <ul style="list-style-type: none"><li>• if epoch &lt; current epoch: send no votes</li><li>• if vote undelegated: send self vote</li><li>• if candidate: send self vote</li><li>• if delegate: send all votes</li></ul> <p>Vote: (epoch <b>e</b>, quorum <b>q</b>, term <b>t</b>, votes <b>v</b>)</p>	<p>Delegations are only valid for the next epoch change. If enough delegates have failed that the epoch change cannot be made, a “nuclear” option resets delegates.</p> <p>Triggered by a nuclear timeout <math>\gg</math> root election timeout to ensure root leader is dead and delegates can’t establish leader.</p> <ul style="list-style-type: none"><li>• Increment epoch beyond vote delegation limit, resetting all delegations.</li><li>• Conduct new root election/epoch change with all available replicas.</li><li>• Update health of all failed nodes and reconfigure epoch.</li></ul>
Epoch Decisions	Epoch Changes	Fuzzy Transitions
	<p>Initiated by request, reconfiguration, localization, quiescence procedures.</p> <p><b>Root Leader</b></p> <ul style="list-style-type: none"><li>• Monotonically increase epoch number, Define members, assign initial leaders.</li><li>• Initiate delegated vote on epoch-change.</li><li>• On commit, begin fuzzy transition.</li></ul> <p><b>Subquorum Replicas</b></p> <ul style="list-style-type: none"><li>• Write tombstone into current log.</li><li>• Finalize commit for accesses prior to the tombstone record, forward new requests.</li><li>• On tombstone commit: truncate and archive log, join new subquorum configuration.</li></ul>	<p><b>Initiating:</b> leader of subquorum in e-1</p> <p><b>Remote:</b> leader of subquorum in e</p> <ul style="list-style-type: none"><li>• Initiating sends last committed command for every object required by remote, Null for objects without accesses, and number of outstanding entries.</li><li>• Remote appends last entries and performs batch consensus to bring subquorum to the Same state.</li><li>• On remote commit, reports to root leader and begins accepting new accesses.</li></ul> <p><b>Note:</b> background anti-entropy optimizes handoff process by reducing data volume.</p>
Operations	Consensus and Accesses	Remote Accesses
	<p>Clients are forwarded to the subquorum leader with responsibility for requested object(s).</p> <ul style="list-style-type: none"><li>• <b>Read(o):</b> Leader responds with last committed entry; marks response if uncommitted entry for object exists. Adds read access to log but does not begin consensus (aggregates reads with writes).</li><li>• <b>Write(o):</b> Leader increments objects version number and creates a corresponding log entry. Sends consensus request and responds to client when the entry is committed.</li></ul>	<p>In a multi-object transaction, remote accesses serialize inter-quorum access.</p> <p><b>Initiating:</b> append entries in log and send remote access request to remote leader.</p> <p><b>Remote:</b> create sub-epoch to demarcate remote access, add entry and respond to initiating replica when committed.</p> <p><b>Initiating:</b> on remote commit, create local sub-epoch, and commit entries appended to logs.</p>

Figure 3.2: A condensed summary of the hierarchical consensus protocol. Operations are described in a top-to-bottom fashion where the top level is root quorum operations, the bottom is subquorum operations, and the middle is transition and intersection.

uses Raft as a building block, so we describe the relevant portions of Raft at a high level, referring the reader to the original paper for complete details. Though we chose to base our protocol on Raft, a similar approach could be used to modify Paxos into a hierarchical structure.

Consensus protocols typically have two phases: leader *election* and operations *commit*<sup>1</sup>. Raft is a strong-leader consensus protocol, which allows the election phase to be elided while a leader remains available. The protocol requires only a single communication round to commit an operation in the common case. Raft uses timeouts to trigger phase changes and provide fault tolerance. Crucially, it relies on timeouts only to provide progress, not safety. New elections occur when another replica in the quorum times out waiting for communication from the leader. Such a replica increments its *term* until it is greater than the existing leader, and announce its candidacy. Other replicas vote for the candidate if they have not seen a competing candidate with a larger term. During regular operation, clients send requests to the leader, which broadcasts **AppendEntries** messages carrying operations to all replicas. An operation is *committed* and can be executed when the leader receives acknowledgments of the **AppendEntries** message from more than half the replicas (including itself).

We describe differences in our Raft implementation from the canonical implementation in Chapter 5

Throughout the rest of this chapter we use the term *root quorum* to refer to the upper, namespace-mapping and configuration-management tier of HC, and

---

<sup>1</sup>Election and commit phases correspond to PROPOSE and ACCEPT phases in Paxos

*subquorum* to describe a group of replicas (called *peers*) participating in consensus decisions for a section of the namespace. The root quorum shepherds subquorums through *epochs*, each with potentially different mappings of the namespace and replicas to subquorums. An epoch corresponds to a single commit phase of the root quorum. We use the term Raft only when describing details particular to our current use of Raft as the underlying consensus algorithm. We refer to the two phases of the base consensus protocol as the *election phase* and the *commit phase*. We use the term *vote* as a general term to describe positive responses in either phase. Epoch  $x$  is denoted  $e_x$ . Subquorum  $i$  of epoch  $e_x$  is represented as  $q_{i,x}$ , or just  $q_i$  when the epoch is obvious.  $t_a$  represents a specific *tag*, or disjoint subset of the namespace.

We assume faults are fail-stop [49] rather than Byzantine [58]. We do not assume that either replica hosts or networks are homogeneous, nor do we assume freedom from partitions and other network faults.

### 3.2.1 Root Consensus

Hierarchical consensus is a leader-oriented protocol that organizes replicas into two tiers of quorums, each responsible for fundamentally different decisions (Figure 3.3). The lower tier consists of multiple independent subquorums, each committing operations to local shared logs. The upper, *root quorum*, consists of subquorum peers, usually their leaders, delegated to represent the subquorum in root elections and commits. Hierarchical consensus’s main function is to export a linearizable abstraction of shared accesses to some underlying substrate, such as a

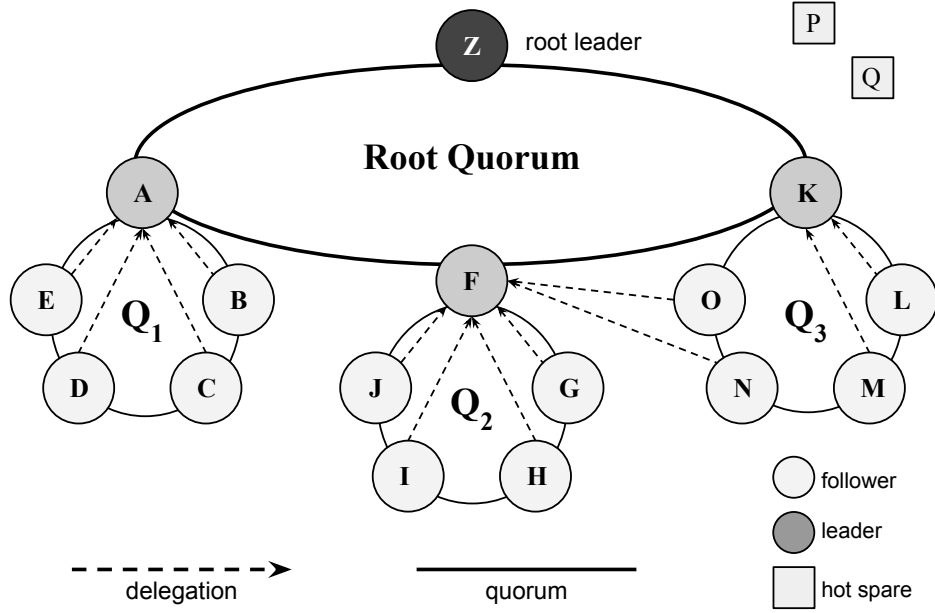


Figure 3.3:

distributed object store or file system. We assume that nodes hosting object stores, applications, and HC are frequently co-located across the wide area.

The root quorum’s primary responsibilities are mapping namespaces and replicas to individual subquorums. Each such map defines a distinct epoch,  $e_i$ , a monotonically increasing representation of the term of  $q_{i,e}$ . The root quorum is effectively a consensus group consisting of subquorum leaders. Somewhat like subquorums, the effective membership of the root quorum is not defined by the quorum itself, but in this case by leader election or peer delegations in the lower tier.

The root quorum partitions the namespace across multiple subquorums, each with a disjoint portion as its scope. The intent of subquorum localization is ensure that the *domain* of a client, the portion of the namespace it accesses, is entirely

within the scope of a local, or nearby, subquorum. If true across the entire system, each client interacts with only one subquorum, and subquorums do not interact at all during execution of a single epoch. This *siloing* of client accesses simplifies implementation of strong consistency guarantees and allows better performance.

### 3.2.2 Delegation

Fault tolerance scales with increasing system size. The root quorum’s membership is, at least logically, the set of all system replicas, at all times. However, running consensus elections across large systems is inefficient in the best of cases, and prohibitively slow in a geo-replicated environment. Root quorum decision-making is kept tractable by having replicas *delegate* their votes, usually to their leaders, for a finite duration. With leader delegation, the root membership effectively consists of the set of subquorum leaders. Each leader votes with a count describing its own and peer votes from its subquorum.

Consider an alternative leader-based approach where root quorum membership is defined as the current set of subquorum leaders. Both delegation and the leader approach have clear advantages in performance and flexibility over direct votes of the entire system. However, the leader approach dramatically decreases fault tolerance. Furthermore, the root quorum becomes unstable in the leader approach as its membership changes during partitions or subquorum elections. These changes would require heavyweight *joint consensus* decisions in the root quorum for correctness in Raft-like protocols [48].

With delegation, however, root quorum membership is always the entire system and remains unchanged over subquorum re-configuration. Delegation is essentially a way to optimistically shortcut contacting every replica for each decision. Subquorum repartitioning merely implies that a given replica’s vote might need to be delegated to a different leader.

Delegation does add one complication: the root quorum leader must know all vote delegations to request votes when committing epoch changes. We deal with this issue, as well as the requirement for a nuclear option (Section 3.4.3), by simplifying our protocol. Instead of sending vote requests just to subquorum leaders, **the root quorum leader sends vote requests to all system replicas**. This is true even for *hot spares*, which are not currently in any subquorum.

This is correct because vote requests now reach all replicas, and because replicas whose votes have been delegated merely ignore the request. We argue that it is also efficient, as a commit’s efficiency depends only on receipt of a majority of the votes. Large consensus groups are generally slow (see Section 3.5) not just because of communication latency, but because large groups in a heterogeneous setting are more likely to include replicas on very slow hosts or networks. In the usual case for our protocol, the root leader still only needs to wait for votes from the subquorum leaders. Leaders are generally those that respond more quickly to timeouts, so the speed of root quorum operations is unchanged.

### 3.2.3 Epoch Transitions

An epoch change is initiated by the leader in response to one of several events, including:

- a namespace repartition request from a subquorum leader
- notification of join requests by new replicas
- notification of failed replicas
- changing network conditions that suggest re-assignment of replicas

The root leader transitions to a new epoch through the normal commit phase in the root quorum. The command proposed by the leader is an enumeration of the new subquorum partition, namespace partition, and assignment of namespace portions to specific subquorums. The announcement may also include initial leaders for each subquorum, with the usual rules for leader election applying otherwise, or if the assigned leader is unresponsive. Upon commit, the operation serves as an *announcement* to subquorum leaders. Subquorum leaders repeat the announcement locally, disseminating full knowledge of the new system configuration, and eventually transition to the new epoch by committing an **epoch-change** operation locally.

The epoch change is lightweight for subquorums that are not directly affected by the underlying re-configuration. If a subquorum is being changed or dissolved, however, the *epoch-change* commitment becomes a tombstone written to the logs of all local replicas. No further operations will be committed by that version of the



subgroup, and the local shared log is archived and then truncated. Truncation is necessary to guarantee a consistent view of the log within a subquorum, as peers may have been part of different subquorums, and thus have different logs, during the last epoch. Replicas then begin participating in their new subquorum instantiation. In the common case where a subquorum's membership remains unchanged across the transition, an **epoch-change** may still require additional mechanism because of changes in namespace responsibility.

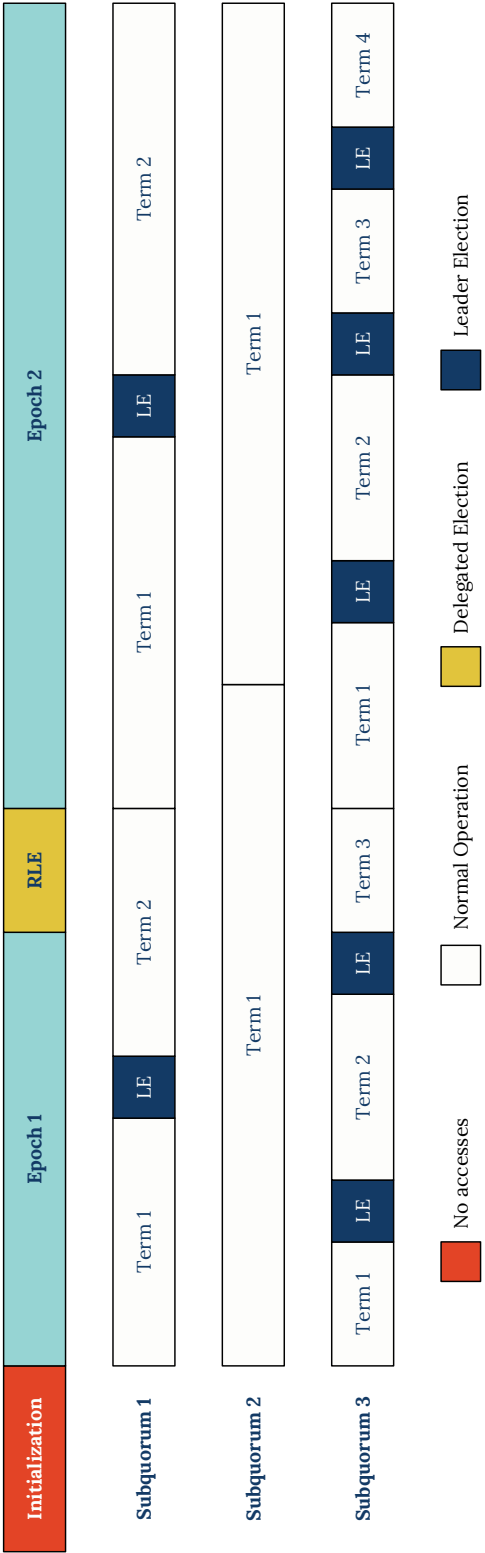


Figure 3.4:

### 3.2.4 Fuzzy Handshakes

Epoch handshakes are required whenever the namespace-to-subquorum mapping changes across an epoch boundary. HC separates epoch transition announcements in the root quorum from implementation in subquorums. Epoch transitions are termed *fuzzy* because subquorums need not all transition synchronously. There are many reasons why a subquorum might be slow. Communication delays and partitions might delay notification. Temporary failures might block local commits. A subquorum might also delay transitioning to allow a local burst of activity to cease such as currently running transactions <sup>2</sup>. Safety is guaranteed by tracking subquorum dependencies across the epoch boundary.

Figure 3.5 shows an epoch transition where the scopes of  $q_i$ ,  $q_j$ , and  $q_k$  change across the transition as follows:

Fix the alignment of the below equations:

$$q_{i,x-1} = t_a, t_b \quad \longrightarrow \quad q_{i,x} = t_a \quad (3.1)$$

$$q_{j,x-1} = t_c, t_d \quad \longrightarrow \quad q_{j,x} = t_c, t_d, t_f \quad (3.2)$$

$$q_{k,x-1} = t_e, t_f \quad \longrightarrow \quad q_{k,x} = t_d, t_e \quad (3.3)$$

All three subquorums learn of the epoch change at the same time, but become ready with varying delays. These delays could be because of network lags or ongoing local activity. Subquorum  $q_i$  gains no new tags across the transition and moves

---

<sup>2</sup>The HC implementation discussed in this chapter does not currently support transactions.

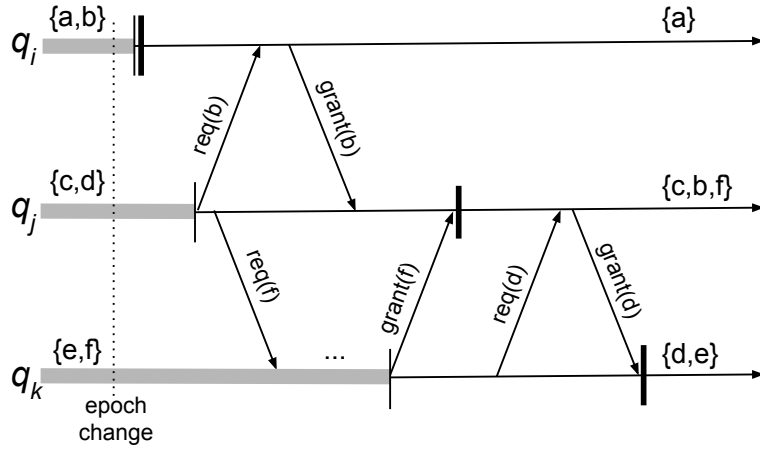


Figure 3.5: Readiness to transition to the new epoch is marked by a thin vertical bar; actual transition is the thick vertical bar. Thick gray lines indicate operation in the previous epoch. Subquorum  $q_j$  transitions from tag  $c, d$  to  $c, b, f$ , but begins only after receiving version information from previous owners of those tags. The request to  $q_k$  is only answered once  $q_k$  is ready to transition as well.

immediately to the new epoch. Subquorum  $q_j$ 's readiness is slower, but then it sends requests to the owners of both the new tags it acquires in the new epoch. Though  $q_i$  responds immediately,  $q_k$  delays its response until locally operations conclude. Once both handshakes are received,  $q_j$  moves into the new epoch, and  $q_k$  later follows suit.

These bilateral handshakes allow an epoch change to be implemented incrementally, eliminating the need for lockstep synchronization across the entire system. This flexibility is key to coping with partitions and varying connectivity in the wide area. However, this piecewise transition, in combination with subquorum re-definition and configuration at epoch changes, also means that individual replicas *may be part of multiple subquorums at a time*.

This overlap is possible because replicas may be mapped to distinct subgroups from one epoch to the next. Consider  $q_k$  in Figure 3.5 again. Assume the epochs shown are  $e_x$  and  $e_{x+1}$ . A single replica,  $r_a$ , may be remapped from subquorum  $q_{k,x}$  to subquorum  $q_{i,x+1}$  across the transition. Subquorum  $q_{k,x}$  is late to transition, but  $q_{i,x+1}$  begins the new epoch almost immediately. Requiring  $r_a$  to participate in a single subquorum at a time would potentially delay  $q_{i,x+1}$ 's transition and impose artificial synchronicity constraints on the system. One of the many changes we made in the base Raft protocol is to allow a replica to have multiple distinct shared logs. Smaller changes concern the mapping of requests and responses to the appropriate consensus group.

### 3.2.5 Subquorum Consensus

this section is bad

Each subquorum,  $q_i$ , elects a leader to coordinate local decisions. Fault tolerance of the subquorum is maintained in the usual way, detecting leader failures and electing new leaders from the peers. Subquorums do not, however, ever change system membership on their own. Subquorum membership is always defined in the root quorum.

Subquorum consensus is used to commit object writes. Reads are not committed by default, but are always served by the leader of the appropriate subquorum. Namespace assignments in HC result in the object space being partitioned (or sharded) across distinct subquorums. The mapping of the shared object namespace to individual subquorums is the *tagset*, or the tagset partition. An individual *tag* defines a disjoint subset of the object space.

As writes are committed through HC, the shared logs provide a complete version history of all distributed objects. Subquorum leaders use in-core caches to provide fast access to recently accessed objects in the local subquorums's tag. Replicas perform background anti-entropy [11, 59, 60], disseminating log updates a user-defined number of times across the system.

The most complex portion of the HC protocol is in handling data-related issues at epoch transitions. Transitions may cause tags to be transferred from one subquorum to another, forcing the new leader to load state remotely to serve object requests. Transitions handshakes are augmented in three ways. First, an replica

can demand-fetch an object version from any other system replica. Second, epoch handoffs contain enumerations of all current object versions, though not the data itself. Knowing an object’s current version gives the new handler of a tag the ability to demand fetch an object that is not yet present locally. Finally, handshakes start immediate fetches of the in-core version cache from the leader of the tag’s subquorum in the old epoch to the leader in the new.

We do not currently gather the entire shared log onto a single replica because of capacity and flexibility issues. Capacity is limited because our system and applications are expected to be long-lived. Flexibility is a problem because HC, and applications built on HC, gain much of their value from the ability to pivot quickly, whether to deal with changes in the environment or for changing application access patterns. We require handoffs to be as lightweight as possible to preserve this advantage.

### 3.2.6 Client Operations

- Sessions - Connect to closest available replica, redirected to closest available leader.

Client namespace accesses are forwarded to the leader of the subquorum for the appropriate part of the namespace. The underlying Raft semantics ensure that leadership changes do not result in loss of any commits. Hence, individual- or multiple-client accesses to a single subquorum are totally ordered. *Remote accesses*, or client accesses to other than their local subquorum, are transparent to the primary

protocol. However, creation of a single shared log of all system operations requires remote accesses to be logged.

### 3.3 Consistency

Pushing all writes through subquorum commits and serving reads at leaders allows us to guarantee that accesses are linearizable (Lin), which is the strongest non-transactional consistency [61,62]. As a recap, linearizability is a combination of atomicity and timeliness guarantees about accesses to a single object. Both **reads** and **writes** must appear atomic, and also instantaneous at some time between a request and the corresponding response to a client. **Reads** must always return the latest value. This implies that reads return values are consistent *with any observed ordering*, i.e., the ordering is *externalizable* [63].

Linearizability of object accesses can be *composed*. If operations on each object are linearizable, the entire object space is also linearizable. This allows our subquorums to operate independently while providing a globally consistent abstraction.

#### 3.3.1 Globally Consistent Logs

Our default use case is in providing linearizable access to an object store. Though this approach allows us to guarantee all observers will see linearizable results of object accesses in real-time, the system is not able to enumerate a total order, or create a linearizable shared log. Such a linear order would require fine-grained (expensive) coordination across the entire system, or fine-grained clock synchro-



nization [16]. Though many or most distributed applications (objects stores, file systems, etc.) will work directly with HC, shared logs are a useful building block for distributed systems.

HC *can* be used to build a sequentially consistent (SC) shared log as shown in Figure 3.6. Like Lin, SC requires all observers to see a single total ordering. SC differs in that this total ordering does not have to be externalizable. Instead, it merely has to conform to local operation orders and all reads-from dependencies.

write about grid consistency

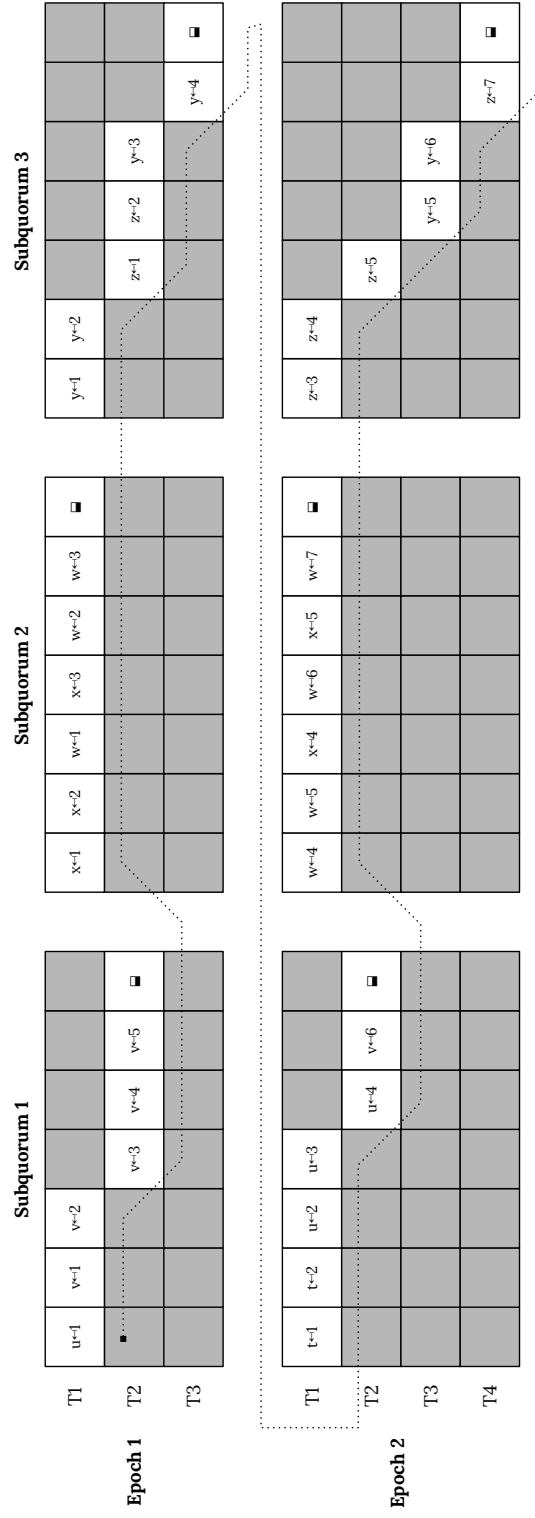


Figure 3.6:

Figure 3.7 shows a system with subquorums  $q_i$  and  $q_j$ , each of which performs a pair of writes. Dotted lines show one possible event ordering for replicas  $q_i$  (responsible for objects  $a$  and  $b$ ), and  $q_j$  ( $c$  and  $d$ ). Without cross-subquorum reads or writes, ordering either subquorum’s operations first creates a SC total ordering:  $q_i \rightarrow q_j$  (“happened-before” [64]) implies  $w_{i,1} \rightarrow w_{i,3} \rightarrow w_{j,1} \rightarrow w_{j,3}$ , for example.

By contrast, the subquorums in Figure 3.8 create additional dependencies by issuing remote writes to other subquorums:  $w_{i,2} \rightarrow w_{j,3}$  and  $w_{j,2} \rightarrow w_{i,3}$ . Each remote write establishes a partial ordering between events of the sender before the sending of the write, and writes by the receiver after the write is received. Similar dependencies result from remote reads.

These dependencies cause the epochs to be split (not shown in picture). The receipt of write  $w_{i,2}$  in  $q_j$  causes  $q_{j,1}$  to be split into  $q_{j,1.1}$  and  $q_{j,1.2}$ . Likewise, the receipt of write  $w_{j,2}$  into  $q_i$  causes  $q_i$  to be split into  $q_{i,1.1}$  and  $q_{i,1.2}$ . Any topological sort of the subepochs that respects these orderings, such as  $q_{i,1.1} \rightarrow q_{j,1.1} \rightarrow q_{j,1.2} \rightarrow q_{i,1.2}$ , results in a valid SC ordering.

Presenting a sequentially consistent global log across the entire system, then, only requires tracking these inter-subquorum data accesses, and then performing an  $\mathcal{O}(n)$  merge of the subepochs.

By definition, this log’s ordering respects any externally visible ordering of cross-subquorum accesses (accesses visible to the system). However, the log does not necessarily order other accesses according to external visibility. The resulting shared log could not be mined to find causal relationships between accesses through external communication paths unknown to the system.

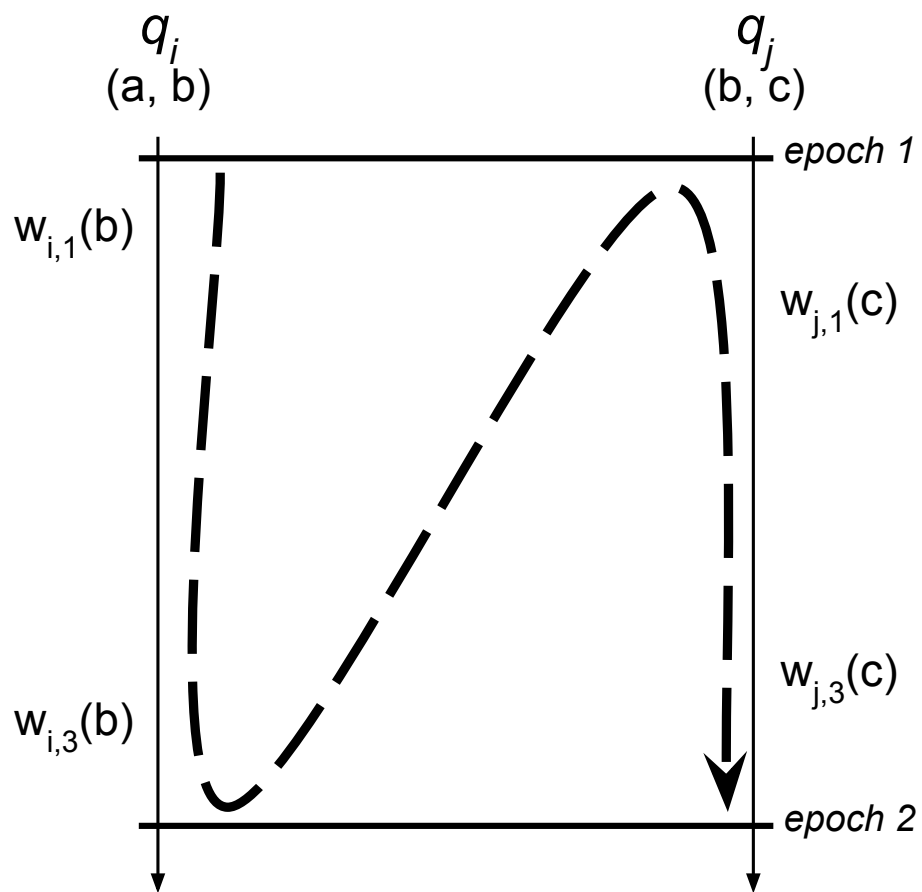


Figure 3.7: Default ordering:  $w_{i,1} \rightarrow w_{i,3} \rightarrow w_{j,1} \rightarrow w_{j,3}$

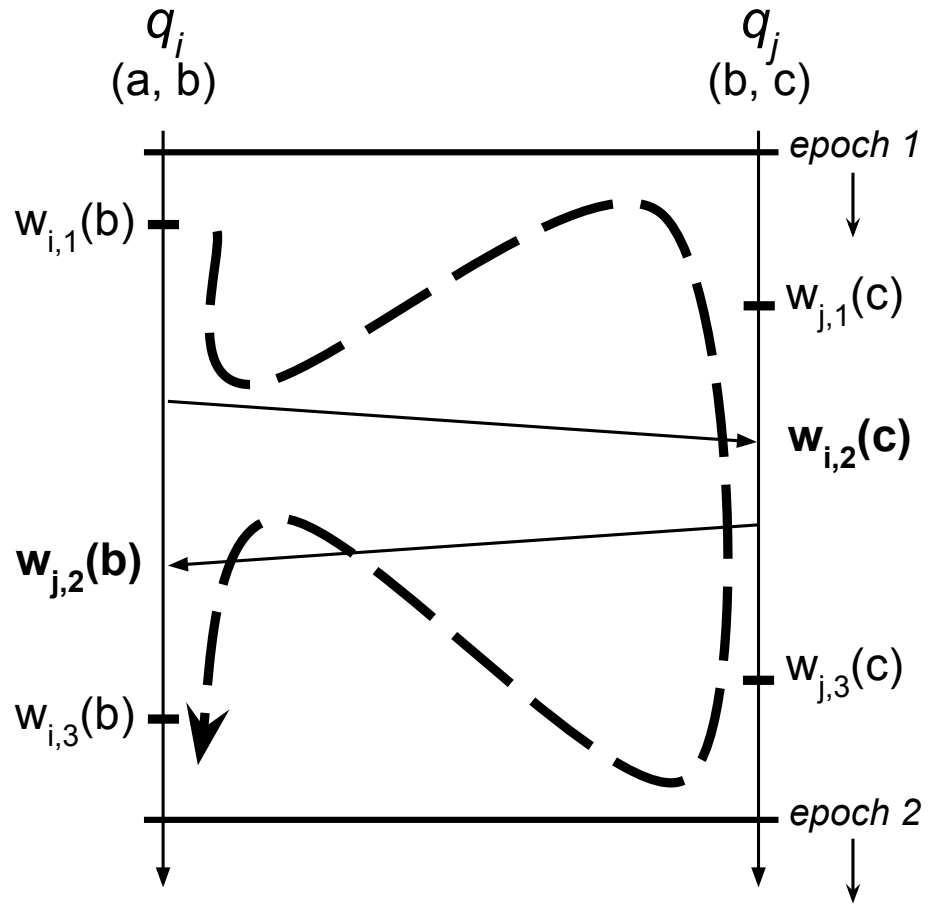


Figure 3.8: Remote writes add additional ordering constraints:  $w_{i,1} \rightarrow w_{i,2} \rightarrow w_{j,3}$ , and  $w_{j,1} \rightarrow w_{j,2} \rightarrow w_{i,3}$

For example, assume that log events are published posts, and that one user claimed plagiarism. The accused would not be able to prove that his post came first unless there were some causal chain of posts and references visible to the protocol.

### 3.4 Fault Tolerance

We assert that consensus at the leaf replicas is correct and safe because decisions are implemented using well-known leader-oriented consensus approaches. Hierarchical consensus therefore has to demonstrate linearizable correctness and safety between subquorums for a single epoch and between epochs. Briefly, linearizability requires external observers to view operations to objects as instantaneous events. Within an epoch, subquorum leaders serially order local accesses, thereby guaranteeing linearizability for all replicas in that quorum.

Epoch transitions raise the possibility of portions of the namespace being re-assigned from one subquorum to another, with each subquorum making the transition independently. Correctness is guaranteed by an invariant requiring subquorums to delay serving newly acquired portions of the namespace until after completing all appropriate handshakes.

Table 3.1: Failure categories: Peer failure is detected by missed heartbeat messages. The rest are triggered by the appropriate election timeout.

Failure Type	Response
subquorum peer	request replica repartition from root quorum
subquorum leader	local election, request replacement from root quorum
root leader	root election (with delegations)
majority of majority of subquorums	(nuclear option) root election after delegations timed out

### 3.4.1 Failures

During failure-free execution, the root quorum partitions the system into disjoint subquorums, assigns *subquorum leaders*, and assigns partitions of the tag-space to subquorums. Each subquorum coordinates and responds to accesses for objects in its assigned tag-space. We define the system’s *safety* property as guaranteeing that non-linearizable (or non-sequentially-consistent, see Section 3.3.1) event orderings can never be observed. We define the system’s *progress* property as the system having enough live replicas to commit votes or operations in the root quorum.

The system can suffer several types of failures, as shown in Table 3.1. Failures of subquorum and root quorum leaders are handled through the normal consensus mechanisms. Failures of subquorum peers are handled by the local leader petitioning the root quorum to re-configure the subquorum in the next epoch. Failure of a root quorum peer is the failure of subquorum leader, which is handled as above. Root quorum heartbeats help inform other replicas of leadership changes, potentially necessary when individual subquorums break down.

describe assassination

HC’s structure means that some faults are more important than others. Proper operation of the root quorum requires the majority of replicas in the majority of subquorums to be non-faulty. Given a system with  $2m+1$  subquorums, each of  $2n+1$  replicas, the entire system’s progress can be halted with as few as  $(m+1)(n+1)$  well-chosen failures. Therefore, in worst case, the system can only tolerate:  $f_{worst} = mn + m + n$  failures and still make progress. At maximum, HC’s basic protocol



can tolerate up to:  $f_{best} = (m + 1) * n + m * (2n + 1) = 3mn + m + n$  failures. As an example, a 25/5 system can tolerate at least 8 and up to 16 failures out of 25 total replicas. A 21/3 system can tolerate at least 7, and a maximum of 12, failures out of 21 total replicas. Individual subquorums might still be able to perform local operations despite an impasse at the global level.

Total subquorum failure can temporarily cause a portion of the namespace to be unserved. However, the root quorum eventually times out and moves into a new epoch with that portion assigned to another subquorum.

### 3.4.2 Obligations Timeout

write this section

### 3.4.3 The Nuclear Option

fix this section

Singleton consensus protocols, including Raft, can tolerate just under half of the entire system failing. As described above, HC's structure makes it more vulnerable to clustered failures. Therefore we define a *nuclear option*, which uses direct consensus decision among all system replicas to tolerate any  $f$  replicas failing out of  $2f + 1$  total replicas in the system.

A nuclear vote is triggered by the failure of a root leader election. A *nuclear candidate* increment's its term for the root quorum and broadcasts a request for votes to all system replicas. The key difficulty is in preventing delegated votes

and nuclear votes from reaching conflicting decisions. Such situations might occur when temporarily unavailable subquorum leaders regain connectivity and allow a wedged root quorum to unblock. Meanwhile, a nuclear vote might be concurrently underway.

Replica delegations are defined as intervals over specific slots. Using local subquorum slots would fall prey to the above problem, so we define delegations as a small number (often one) of root slots, which usually correspond to distinct epochs. During failure-free operation, peers delegate to their leaders and are all represented in the next root election or commit. Peers then renew their delegations to their leaders by appending them to the next local commit reply. This approach works for replicas that change subquorums over an epoch boundary, and even allows peers to delegate their votes to arbitrary other peers in the system (see replicas  $r_N$  and  $r_O$  in Figure 3.3).

This approach is simple and correct, but deals poorly with leader turnovers in the subquorum. Consider a subquorum where all peers have delegated votes to their leader for the next root slot. If that leader fails, none of the peers will be represented. We finesse this issue by re-defining such delegations to count root elections, root commits, *and* root heartbeats. The latter means that local peers will regain their votes for the next root quorum action if it happens after to the next heartbeat.

Consider the worst-case failure situation discussed in Section 3.4: a majority of the majority of subquorums have failed. None of the failed subquorum leaders can be replaced, as none of those subquorums have enough local peers.

The first response is initiated when a replica holding delegations (or its own vote) times out waiting for the root heartbeat. That replica increments its own root term, adopts the prior system configuration as its own, and becomes a root candidate. This candidacy fails, as a majority of subquorum leaders, with all of their delegated votes, are gone. Progress is not made until delegations time out. In our default case where a delegation is for a single root event, this happens after the first root election failure.

At the next timeout, any replica might become a candidate because delegations have lapsed (under our default assumptions above). Such a *nuclear* candidate increments its root term and sends candidate requests to all system replicas, succeeding if it gathers a majority across all live replicas.

The first candidacy assumed the prior system configuration in its candidacy announcement. This configuration is no longer appropriate unless some of the “failed” replicas quickly regain connectivity. Before the replica announces its candidacy for a second time, however, many of the replica replies have timed out. The candidate alters its second proposed configuration by recasting all such replicas as hot spares and potentially reducing the number and size of the subgroups. Subsequent epoch changes might re-integrate the new hot spares if the replicas regain connectivity.

### 3.5 Performance Evaluation

HC was designed to adapt both to dynamic workloads as well as variable network conditions. We therefore evaluate HC in two distinct environments: a

homogeneous data center and a heterogeneous real-world network. The homogeneous cluster is hosted on Amazon EC2 and includes 26 “t2.medium” instances: dual-core virtual machines running in a single VPC with inter-machine latencies of  $\lambda_\mu = 0.399ms$  and  $\lambda_\sigma = 0.216ms$ . These machines are cost effective and, though lightweight, are easy to scale to large cluster sizes as workload increases. Experiments are set up such that each instance runs a single replica process and multiple client processes.

The heterogeneous cluster (UMD) consists of several local machines distributed across a wide area, with inter-machine latencies ranging from  $\lambda_\mu = 2.527ms$ ,  $\lambda_\sigma = 1.147ms$  to  $\lambda_\mu = 34.651ms$ ,  $\lambda_\sigma = 37.915ms$ . Machines in this network are a variety of dual and quad core desktop servers that are solely dedicated to running these benchmarks. Experiments on these machines are set up so that each instance runs multiple replica and client processes co-located on the same host. In this environment, localization is critical both for performance but also to ensure that the protocol can elect and maintain consensus leadership. The variability of this network also poses challenges that HC is uniquely suited to handle via root quorum-guided adaptation. We explore two distinct scenarios – sawtooth and repartitioning – using this cluster; all other experiments were run on the EC2 cluster.

HC is partially motivated by the need to scale strong consistency to large cluster sizes. We based our work on the assumption that consensus performance decreases as the quorum size increases, which we confirm empirically in Figure 3.9. This figure shows the maximum throughput against system size for a variety of workloads, up to 120 concurrent clients. A workload consists of one or more clients

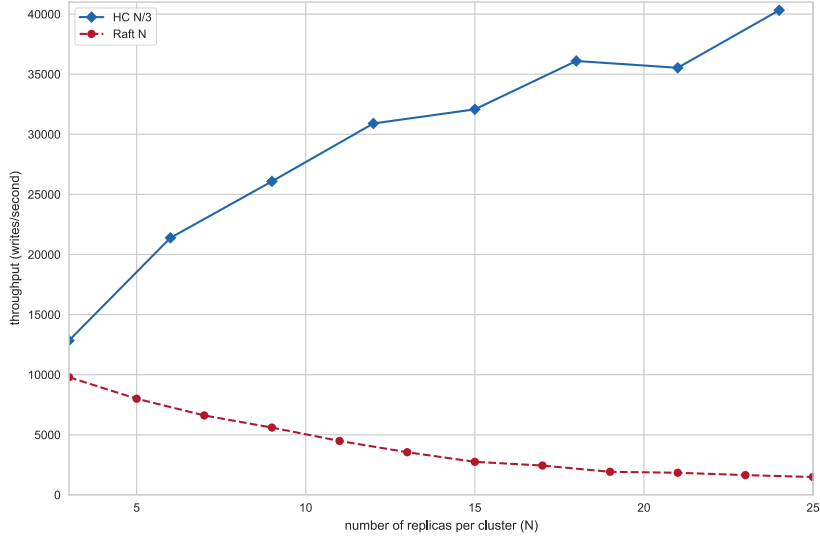


Figure 3.9: Mean throughput of workloads of up to 120 concurrent clients

continuously sending writes of a specific object or objects to the cluster without pause.

Standard consensus algorithms, Raft in particular, scale poorly with uniformly decreasing throughput as nodes are added to the cluster. Commit latency increases with quorum size as the system has to wait for more responses from peers, thereby decreasing overall throughput. Figures 3.9 and 3.10 clearly show the multiplicative advantage of HC’s hierarchical structure, though HC does not scale linearly as we had expected.

There are at least two factors currently limiting the HC throughput shown here. First, the HC subquorums for the larger system sizes are not saturated. A single 3-node subquorum saturates at around 25 clients and this experiment has only

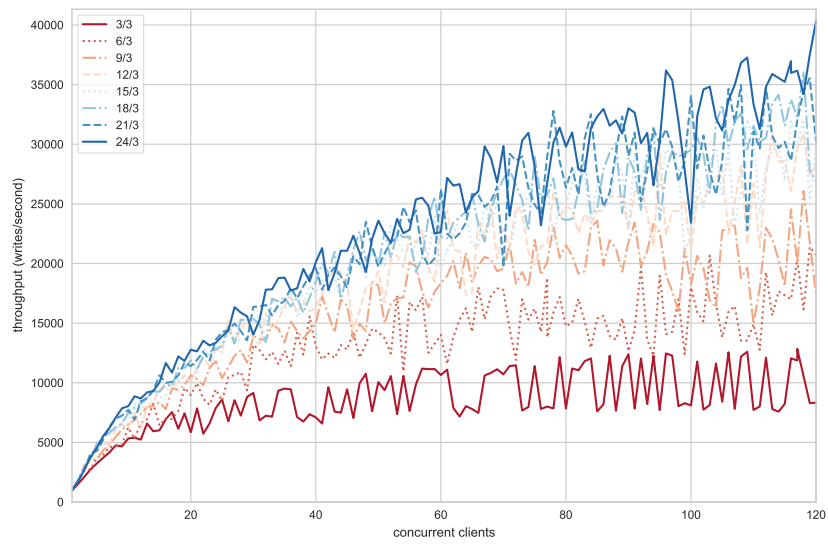


Figure 3.10: Performance of distributed consensus with an increasing workload of concurrent clients. Performance is measured by throughput, the number of writes committed per second.

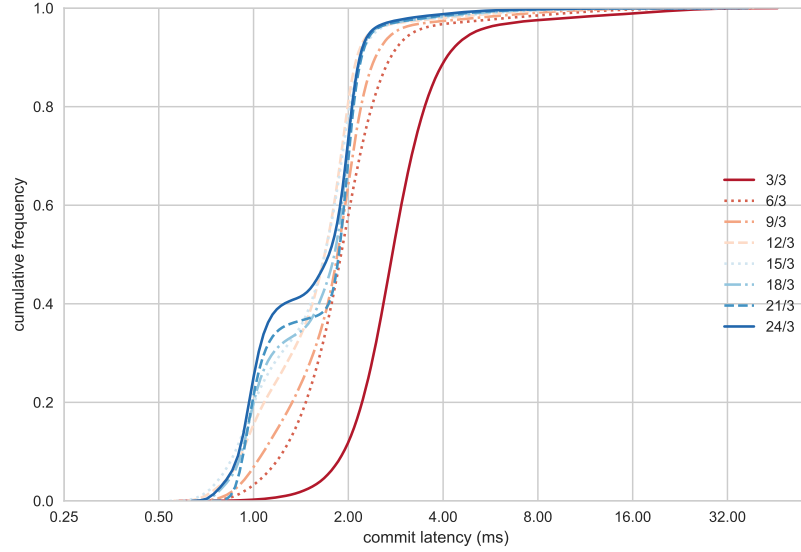


Figure 3.11:

about 15 clients per subquorum for the largest cluster size. We ran experiments with 600 clients, saturating all subquorums even in the 24-node case. This throughput peaked at slightly over 50,000 committed writes per second, better but still lower than the linear scaling we had expected.

We think the reason for this ceiling is hinted at by Figure 3.10. This figure shows increasingly larger variability with increasing system sizes. A more thorough examination of the data shows widely varying performance across individual subquorums in the larger configurations. We suspect that the cause is either VM misconfiguration or misbehavior. We are adding more instrumentation to diagnose the problem.

The effect of saturation is also demonstrated in Figure 3.11, which shows cumu-

lative latency distributions for different system sizes holding the workload (number of concurrent clients) constant. The fastest (24/3) shows nearly 80% of client write requests being serviced in under 2 msec. Larger system sizes are faster because the smaller systems suffer from contention (25 clients can saturate a single subquorum). Because throughput is directly related to commit latency, throughput variability can be mitigated by adding additional subquorums to balance load.

Besides pure performance and scaling, HC is also motivated by the need to adapt to varying environmental conditions. In the next set of experiments, we explore two common runtime scenarios that motivate adaptation: shifting client workloads and failures. We show that HC is able to adapt and recover with little loss in performance. These scenarios are shown in Figures 3.12 and 3.13 as throughput over time, where vertical dotted lines indicate an epoch change.

The first scenario, described by the time series in Figure 3.12 shows an HC 3-replica configuration moving through two epoch changes. Each epoch change is triggered by the need to localize tags accessed by clients to nearby subquorums. The scenario shown starts with all clients co-located with the subquorum serving the tag they are accessing. However, clients incrementally change their access patterns first to a tag located on one remote subquorum, and then to the tag owned by the other. In both cases, the root quorum adapts the system by repartitioning the tagspace such that the tag defining their current focus is served by the co-located subquorum.

Finally, Figure 3.13 shows a 3-subquorum configuration where one entire subquorum becomes partitioned from the others. After a timeout, the root uses an epoch change to re-allocate the tag of the partitioned subquorum over the two



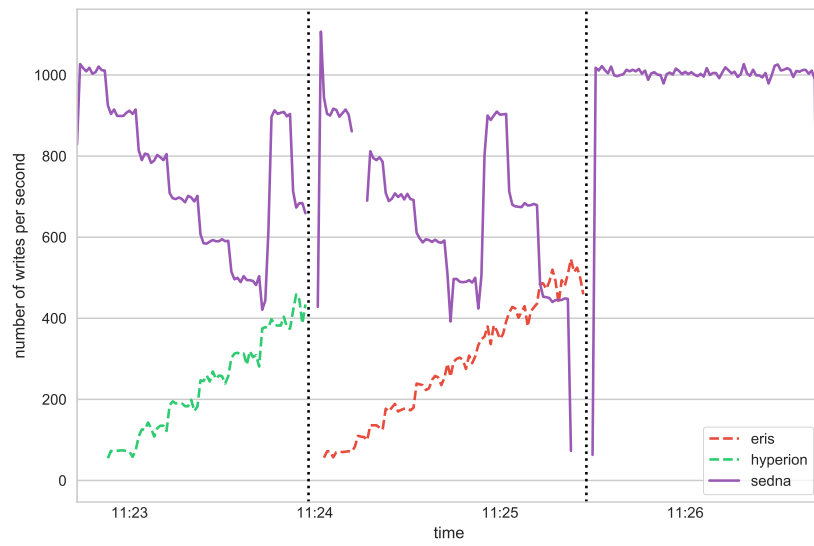


Figure 3.12: 9/3 system adapting to changing client access patterns by repartitioning the tag space so that clients are co-located with subquorums that serve tags they need.

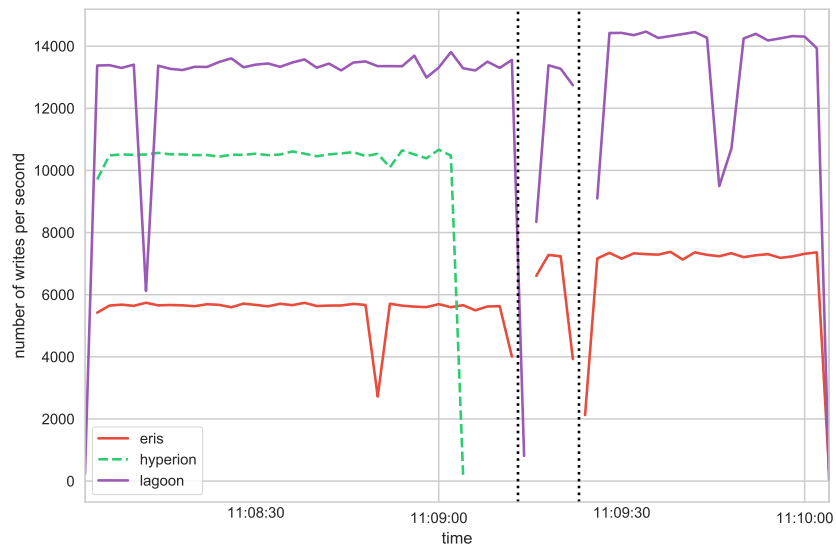


Figure 3.13: 9/3 System that adapts to failure (partition) of entire subquorum. After timeout, the root quorum re-partitions the tag allocated to the failed subquorum among the other two subquorums.

remaining subquorums. The partitioned subquorum eventually has an *obligation timeout*, after which the root quorum is not obliged to leave the tag with the current subquorum. The tag may then be re-assigned to any other subquorum. Timeouts are structured such that by the time an obligation timeout fires, the root quorum has already re-mapped that subquorum's tag to other subquorums. As a result, the system is able to recover from the partition as fast as possible. Note that in this figure, the repartition occurs through two epoch changes, the first allocating part of the tagspace to the first subquorum, and the second allocating the rest of the tag to the other. Gaps in the graph are periods where the subquorums are electing local leaders. This may be optimized by having leadership assigned or maintained through root consensus.

### 3.6 Conclusion

Most consensus algorithms have their roots in the Paxos algorithm, originally described in parliamentary terms. The metaphor of government still applies well as we look at the evolution of distributed coordination as systems have grown to include large numbers of processes and geographies. Systems that use a dedicated leader are easy to reason about and implement, however, like chess, if the leader goes down the system cannot make any progress. Simple democracies for small groups solve this problem but do not scale, and as the system grows, it fragments into tribes. Inspired by modern governments, we have proposed a representative system of consensus, hierarchical consensus, such that replicas elect leaders to participate

in a root quorum that makes decisions about the global state of the system. Local decision making, the kind that effects only a subset of clients and objects is handled locally by subquorums as efficiently as possible. The result is a hierarchy of decision making that takes advantage of hierarchies that already exist in applications.

Hierarchical Consensus is an implementation and extension of Vertical Paxos. Like Vertical Paxos, HC reasons about consistency across all objects by identifying commands with a grid ordering (rather than a log ordering) and is reconfigurable to adapt to dynamic environments that exist in geo-replicated systems. Adaptability allows HC to exploit locality of access, allowing for high performance coordination, even with replication across the wide area. HC extends Vertical Paxos to ensure that intersections exist between the subquorums and the root quorum to ensure coordination exists between subquorums and to ensure that the system operates as a coordinated whole. To scale the consensus protocol of the root quorum, we propose a novel approach, delegation, to ensure that all replicas participate in consensus but limit the number and frequency of messages required to achieve majority. Finally, we generalized HC from primary-backup replication to describe more general online replication required by distributed databases and file systems.

In the next chapter we'll explore a hybrid consistency model implemented by federating replicas that participate in different consistency protocols. In a planetary scale network, HC provides the strong consistency backbone of the federated model, increasing the overall consistency of the system by making coordinating decisions at a high level, and allowing high availability replicas in the fog operate independently where necessary.

## Chapter 4: Federated Consistency

Hybrid consistency model

### 4.1 Overview

### 4.2 Eventual Consistency

Read/Write Quorums

Anti-Entropy Sessions and Synchronization

Policies: latest writer wins

Bilateral anti-entropy

#### 4.2.1 Consistency Failures

Forks

Stale Reads

## 4.3 Integration

### 4.3.1 Communication Integration

### 4.3.2 Consistency Integration

- Forte Number

## 4.4 Performance Evaluation

Communication Topology Inconsistencies due to outages Inconsistency due to system latency

## Chapter 5: System Implementation

Given its grandiose title, it may seem that the engineering behind the development of a planetary scale data storage system would require thousands of man-hours of professional software engineers and a highly structured development process. In fact, this is not necessarily the case for two reasons. First, data systems benefit from an existing global network topology and commercial frameworks for deploying applications. This means that both the foundation and motivation for creating large geo-replicated systems exists, as described earlier. Second, like the internet, complex global systems emerge through the composition of many simpler components following straight forward rules [65]. Instead of architecting a monolithic system, the design process is decomposed to reasoning about the behavior of single processes. Rather than being built, a robust planetary data system evolves from its network environment.

To facilitate system evolution, the consistency models we have described thus far have been *composable* to allow heterogeneous replicas to participate in the same system. However, while composability allows us to reason about consistency expectations, it does not necessarily mean interoperability. In order to ensure reason about system expectations we must outline our assumptions for communication, se-

curity, processing, and data storage. In this chapter, we describe the implementation of the replicas and applications of our experimental system and the assumptions we made.

## 5.1 System Model

A *replica* is an independent process that maintains a portion of the objects stored as well as a *view* of the state of the entire system. Replicas must be able to communicate with one another and may also *serve* requests from clients. A system is composed of multiple communicating replicas and is defined by the behavior the replicas. For example, a totally replicated system is one where each replica stores a complete copy of all objects as in the primary-backup approach [66], whereas a partially replicated system ensures durability such that multiple replicas store the same object but not all replicas store all objects as in the Google File System [67]. At the scale of a multi-region, globally deployed system, we assume that total replication is impractical and primarily consider the partial replication case.

Let's unpack some of the assumptions made by the seemingly simple statements made in the previous paragraph. First, independence means replicas have a shared-nothing architecture [68] and cannot share either memory or disk space. For practical purposes of fault tolerance, we generally assume that there is a one-to-one relationship between a replica and a disk so that a disk failure means only a single replica failure. Second, that each replica must maintain a view of the state of the entire system means both that replicas must be aware of their peers on the network



and that they should know the locations of objects stored on the network. A strict interpretation of this requirement would necessarily make system membership brittle as it would be difficult to add or remove replicas. Alternatively, a centralized interpretation of the view requirement would allow for

Second, the ability to communicate with replicas and serve requests from clients means that the replica must be addressable.

This requirement is seemingly innocuous when taken by itself, however when we also describe a replica as requiring a view of the state of the entire system, it means that a replica must know about the existence of all other replicas in the system. A strict interpretation of having a complete view would necessarily make the system composition brittle, unable to add or remove replicas. Instead we take a less strict view

- networking - actors - event loop - reasons why the above are important -

## 5.2 Applications

### 5.2.1 Distributed Log

### 5.2.2 Key-Value Database

### 5.2.3 File System

## 5.3 Consistency Model

Client-side vs. system-side consistency

Log model of consistency

Continuous consistency scale

Grid consistency model

## 5.4 Raft

Hierarchical Consensus with modified Raft as the underlying consensus protocol exports a linearizable order of accesses to the distributed key-value store. Alia and the HC library are implemented in Golang use gRPC [69] for communication. The system is implemented in **7,924 lines of code**, not including standard libraries or support packages.

An replica implements multiple instantiations of the Raft protocol, which we have modified in several ways. Every replica must run one instantiation of the *root consensus protocol*. Replicas may also run one or more instantiations of the *commit consensus protocol* if they are assigned to a subquorum. Repartition decisions move the system between epochs with a new configuration and tag space, and can only be initiated by messages from peers or monitoring processes. A successful repartition results in a new epoch, tag space, and subquorum topology committed to the root log. **Repartition** messages also serve to notify the network about events that do not require an epoch change, such as the election of a new subquorum leader or bringing a failed node back online.

Each replica implements an event loop that responds to timing events, client requests, and messages from peers. Events may cause the replica to change state,

modify a command log, broadcast messages to peers, modify the key-value store, or respond to a client. Event handlers need to aggressively lock shared state for correctness because Golang and gRPC make extensive use of multi-threading. The balance between correctness and concurrency-driven performance leads to increasing complexity and tighter coupling between components, one that foreshadows extra-process consistency concerns that have been noted in other work [48, 70, 71].

The computing and network environment of a distributed system plays a large role in determining not just the performance of the system, but also its behavior. A simple example is the election timeout parameter of the Raft consensus protocol, which must be much greater than the average time to broadcast and receive responses, and much less than the mean time between failures [48, 72, 73]. If this requirement is not met, leader may be displaced before heartbeat messages arrive, or the system will be unable to recover when a leader fails. As a result, the relationship between timeouts is critically dependent on the mean latency ( $\lambda_\mu$ ) of the network. Howard [74] proposes  $T = \lambda_\mu + 2\lambda_\sigma$  to determine timeouts based on the distribution of observed latencies, sets the heartbeat as  $\frac{T}{2}$ , and the election timeout as the interval  $U(T, 2T)$ . We parameterize our timeouts (Table 5.1) on latency measurements made before we ran our experiments. Monitoring and adapting to network conditions is part of ongoing work.

**Changes to base Raft:** In addition to major changes, such allowing replicas to be part of multiple quorums simultaneously, we also made many smaller changes that had pervasive effects. One change was including the *epoch* number alongside the term in all log entries. The epoch is evaluated for invariants such as whether or

Table 5.1: Parameterized timeouts in our implementation. The *obligation* timeout stops a partitioned subquorum after an extended time without contact to the rest of the system.  $T = 10msec$  for our experiments on Amazon EC2.

Name	Time	Actions
sub heartbeat	1T	sub leader heartbeat
sub leader	2-4T	new sub election
root heartbeat	10T	root leader heartbeat
root election	20-40T	new root election
obligation	50T	root quorum may re-allocate the tag

not a replica can append an entry or if a log is as up to date as a remote log.

Vote delegation requires changes to vote counting. Since our root quorum membership actually consists of the entire system, all replicas are messaged during root events. All replicas reply, though most with a “zero votes” acknowledgment. The root uses observed vote distributions to inform the ordering of future consensus messages (sending requests first to replicas with votes to cast), and uses timeouts to move non-responsive replicas into “hot spares” status.

We allow **AppendEntries** requests in subquorums to aggregate multiple client requests into a single consensus round. Such requests are collected while an outstanding commit round is ongoing, then sent together when that round completes. The root quorum also aggregates all requests within a minimum interval into a single new epoch-change/reconfiguration operation to minimize disruption.

Commits are observed by the leader once a majority of replicas respond positively. Other replicas learn about the commit only on the next message or heartbeat. Root epoch changes and heartbeats are designed to be rare, meaning that epoch

change commits are not seen promptly. We modified the root protocol to inform subquorums of the change by sending an additional heartbeat immediately after it observes a commit.

Replicas may be part of both a subquorum and the root quorum, and across epoch boundaries may be part of multiple subquorums. In principle, a high performance replica may participate in any number of subquorums. We therefore allow replicas to accommodate multiple distinct logs with different access characteristics.

Peers that are either slow or with unsteady connectivity are occasionally left behind at subquorum leader or epoch changes. Root heartbeats containing the current system configuration are broadcast to all replicas and serve to bring them up to date.

Finally, consensus protocols often synchronously write state to disk before responding to remote requests. This allows replicas that merely crash to reboot and rejoin the ongoing computation after recovering state from disk. Otherwise, these replicas need to go through heavyweight leave-and-rejoin handshakes. Our system avoids these synchronous writes by allowing epochs to re-join a subquorum at the next epoch change without any saved state, avoiding these handshakes altogether.

## 5.5 Conclusion

We did not optimize our research for the minimum set of assumptions required to facilitate interoperability between heterogeneous replicas. However, we hope that the assumptions we did make shed light on what is required to achieve the minimum

set of assumptions.

Further research is required to ...

## Chapter 6: Adaptive Consistency

Monitor and Optimize.

Replica placement, object placement

### 6.1 Bandit-Based Approaches

#### 6.1.1 Rewards Function

### 6.2 Access Temperature Approaches

- Expected model of access patterns (daylight).

## Chapter 7: Related Work

Spanner [16] provides global consistency by sharding each tablet across multiple Paxos groups then externalizes their consistency using TrueTime, delaying the commit until a window of uncertainty has passed.

CalvinFS [18, 19] batches transaction operations across the wide area, but still requires paxos to be deployed across the wide area.

Systems that implement many small quorums of coordination [16, 17, 75] avoid the centralization bottleneck and reliability concerns of master-service systems [67, 76] but create silos of independent operation that are not coordinated with respect to each other.

### 7.1 Hierarchical Consensus

Our principle contribution is Hierarchical Consensus, a general technique to compose consensus groups, maintain consistency invariants over large systems, and adapt to changing conditions and application loads. HC is related to the large body of work improving throughput in distributed consensus over the Paxos protocol [42, 47, 77, 78], and on Raft [48, 74]. These approaches focus on fast vs. slow path consensus, eliding phases with dependency resolution, and load balancing.



Our work is also orthogonal in that subquorums and the root quorums can be implemented with different underlying protocols, though the two levels must be integrated quite tightly. Further, HC abstracts reconfiguration away from subquorum consensus, allowing multiple subquorums to move into new configurations and reducing the need for joint consensus [48] and other heavyweight procedures. Finally, its hierarchical nature allows the system to multiplex multiple consensus instances on disjoint partitions of the object space while still maintaining global consistency guarantees.

The global consistency guarantees of HC are in direct contrast to other systems that scale by exploiting multiple consensus instances [12, 16, 17] on a per-object basis. These systems retain the advantage of small quorum sizes but cannot provide system-wide consistency invariants. Another set of systems uses quorum-based decision-making but relaxes consistency guarantees [11, 79, 80]; others provide no way to pivot the entire system to a new configuration [75]. Chain replication [81] and Vertical Paxos [50] are among approaches that control Paxos instances through other consensus decisions. However, HC differs in the deep integration of the two different levels. Whereas these approaches are top down, HC consensus decisions at the root level replace system configuration at the subquorum level, and vice versa.

Possibly the closest system to HC is Scatter [75], which uses an overlay to organize consistent groups into a ring. Neighbors can join, split, and talk amongst themselves. The bottom-up approach potentially allows scaling to many subquorums, but the lack of central control makes it hard to implement global re-maps beyond the reach of local neighbors. HC ties the root quorum and subquorums tightly together,

allowing root quorum decisions to completely reconfigure the running system on the fly either on demand or by detecting changes in network conditions.

We claim very strong consistency across a large distributed system, similar to Spanner [16]. Spanner provides linearizable transactions through use of special hardware and environments, which are used to tightly synchronize clocks in the distributed setting. Spanner therefore relies on a very specific, curated environment. HC targets a wider range of systems that require cost effective scaling in the data center to rich dynamic environments with heterogeneity on all levels.

Finally, shared logs have proven useful in a number of settings from fault tolerance to correctness guarantees. However, keeping such logs consistent in even a single consensus instance has proven difficult [67, 82, 83]. More recent systems are leveraging hardware support to provide fast access to shared logs [18, 19, 53, 55, 56, 84]. To our knowledge, HC is the first work to propose synchronizing shared logs across multiple discrete consensus instances in the wide area.

## Chapter 7: Conclusion

## Appendix A: Formal Specification

Will add formal specification here.

## Bibliography

- [1] John Kubiawicz, David Bindel, Yan Chen, Steven Czerwinski, Patrick Eaton, Dennis Geels, Ramakrishan Gummadi, Sean Rhea, Hakim Weatherspoon, Westley Weimer, and others. Oceanstore: An architecture for global-scale persistent storage. In *ACM Sigplan Notices*, volume 35, pages 190–201.
- [2] Ian Foster and Carl Kesselman. *The Grid 2: Blueprint for a New Computing Infrastructure*. Elsevier.
- [3] Daniel Nurmi, Rich Wolski, Chris Grzegorzczak, Graziano Obertelli, Sunil Soman, Lamia Youseff, and Dmitrii Zagorodnov. The eucalyptus open-source cloud-computing system. In *Cluster Computing and the Grid, 2009. CC-GRID'09. 9th IEEE/ACM International Symposium On*, pages 124–131. IEEE.
- [4] Dirk Merkel. Docker: Lightweight linux containers for consistent development and deployment. 2014(239):2.
- [5] Michael Armbrust, Armando Fox, Rean Griffith, Anthony D. Joseph, Randy Katz, Andy Konwinski, Gunho Lee, David Patterson, Ariel Rabkin, and Ion Stoica. A view of cloud computing. 53(4):50–58.
- [6] Cisco Visual Networking Index. The zettabyte era—trends and analysis.
- [7] Michael J. Casey and Paul Vigna. In blockchain we trust. 15:2018.
- [8] Michael Stonebraker and Joey Hellerstein. What goes around comes around. 4.
- [9] Michael Stonebraker, Daniel Abadi, David J. DeWitt, Sam Madden, Erik Paulson, Andrew Pavlo, and Alexander Rasin. MapReduce and parallel DBMSs: Friends or foes? 53(1):64–71.
- [10] C. Mohan. History repeats itself: Sensible and NonsenSQL aspects of the NoSQL hoopla. In *Proceedings of the 16th International Conference on Extending Database Technology*, pages 11–16. ACM.

- [11] Giuseppe DeCandia, Deniz Hastorun, Madan Jampani, Gunavardhan Kakulapati, Avinash Lakshman, Alex Pilchin, Swaminathan Sivasubramanian, Peter Voshall, and Werner Vogels. Dynamo: Amazon’s highly available key-value store. In *ACM SIGOPS Operating Systems Review*, volume 41, pages 205–220. ACM.
- [12] Fay Chang, Jeffrey Dean, Sanjay Ghemawat, Wilson C. Hsieh, Deborah A. Wallach, Mike Burrows, Tushar Chandra, Andrew Fikes, and Robert E. Gruber. Bigtable: A distributed storage system for structured data. 26(2):4.
- [13] Avinash Lakshman and Prashant Malik. Cassandra: A decentralized structured storage system. 44(2):35–40.
- [14] Ankur Khetrapal and Vinay Ganesh. HBase and Hypertable for large scale distributed storage systems. pages 22–28.
- [15] Jason Baker, Chris Bond, James C. Corbett, J. J. Furman, Andrey Khorlin, James Larson, Jean-Michel Leon, Yawei Li, Alexander Lloyd, and Vadim Yushprakh. Megastore: Providing Scalable, Highly Available Storage for Interactive Services. In *CIDR*, volume 11, pages 223–234.
- [16] James C. Corbett, Jeffrey Dean, Michael Epstein, Andrew Fikes, Christopher Frost, J. J. Furman, Sanjay Ghemawat, Andrey Gubarev, Christopher Heiser, Peter Hochschild, and others. Spanner: Google’s globally distributed database. 31(3):8.
- [17] Tim Kraska, Gene Pang, Michael J. Franklin, Samuel Madden, and Alan Fekete. MDCC: Multi-Data Center Consistency. In *Proceedings of the 8th ACM European Conference on Computer Systems*, pages 113–126. ACM.
- [18] Alexander Thomson, Thaddeus Diamond, Shu-Chun Weng, Kun Ren, Philip Shao, and Daniel J Abadi. Calvin: Fast distributed transactions for partitioned database systems. In *Proceedings of the 2012 ACM SIGMOD International Conference on Management of Data*, pages 1–12. ACM.
- [19] Alexander Thomson and Daniel J. Abadi. CalvinFS: Consistent wan replication and scalable metadata management for distributed file systems. In *13th USENIX Conference on File and Storage Technologies (FAST 15)*, pages 1–14.
- [20] Seth Gilbert and Nancy Lynch. Brewer’s conjecture and the feasibility of consistent, available, partition-tolerant web services. 33(2):51–59.
- [21] Peter Bailis, Aaron Davidson, Alan Fekete, Ali Ghodsi, Joseph M Hellerstein, and Ion Stoica. Highly available transactions: Virtues and limitations. 7(3):181–192.
- [22] Daniel J. Abadi. Consistency tradeoffs in modern distributed database system design: CAP is only part of the story. 45(2):37–42.

- [23] Alexandre Verbitski, Anurag Gupta, Debanjan Saha, Murali Brahmadesam, Kamal Gupta, Raman Mittal, Sailesh Krishnamurthy, Sandor Maurice, Tengiz Kharatishvili, and Xiaofeng Bao. Amazon aurora: Design considerations for high throughput cloud-native relational databases. In *Proceedings of the 2017 ACM International Conference on Management of Data*, pages 1041–1052. ACM.
- [24] stevestein. Scaling out with Azure SQL Database.
- [25] Alain Jobart, Sugu Sougoumarane, Michael Berlin, and Anthony Yeh. Vitess.
- [26] Spencer Kimball, Peter Mattis, and Ben Darnell. CockroachDB.
- [27] Erik Kain. The 'Pokémon GO' Launch Has Been A Complete Disaster [Updated].
- [28] Cloud Datastore.
- [29] Eric A. Brewer. Kubernetes and the path to cloud native. In *Proceedings of the Sixth ACM Symposium on Cloud Computing*, pages 167–167. ACM.
- [30] Luke Stone. Bringing Pokémon GO to life on Google Cloud.
- [31] Makiko Yamazaki. Developer of Nintendo's Pokemon GO aiming for rollout to 200...
- [32] Richard George. Nintendo's President Discusses Region Locking.
- [33] Dropbox — Company Info.
- [34] Slack. Slack About Us.
- [35] WeWork. Global Access.
- [36] Tile About Tile.
- [37] Stella Garber. Lessons Learned From Launching Internationally.
- [38] Desdemona Bandini. RunKeeper Scales to Meet Demand from 24 Million Global Users with New Relic.
- [39] Adam Grossman and Jay LaPorte. Dark Sky Weather App for iOS and Android.
- [40] Matthew Hughes. Signal and Telegram are growing rapidly in countries with corruption problems.
- [41] William J. Bolosky, Dexter Bradshaw, Randolph B. Haagens, Norbert P. Kusters, and Peng Li. Paxos replicated state machines as the basis of a high-performance data store. In *Symposium on Networked Systems Design and Implementation (NSDI)*, pages 141–154.

- [42] Iulian Moraru, David G. Andersen, and Michael Kaminsky. There is more consensus in egalitarian parliaments. In *Proceedings of the Twenty-Fourth ACM Symposium on Operating Systems Principles*, pages 358–372. ACM.
- [43] Yanhua Mao, Flavio Paiva Junqueira, and Keith Marzullo. Mencius: Building efficient replicated state machines for WANs. In *OSDI*, volume 8, pages 369–384.
- [44] Martin Biely, Zoran Milosevic, Nuno Santos, and Andre Schiper. S-paxos: Offloading the leader for high throughput state machine replication. In *Reliable Distributed Systems (SRDS), 2012 IEEE 31st Symposium On*, pages 111–120. IEEE.
- [45] Pierre Sutra and Marc Shapiro. Fast genuine generalized consensus. In *Reliable Distributed Systems (SRDS), 2011 30th IEEE Symposium On*, pages 255–264. IEEE.
- [46] Sebastiano Peluso, Alexandru Turcu, Roberto Palmieri, Giuliano Losa, and Binoy Ravindran. Making fast consensus generally faster. In *Dependable Systems and Networks (DSN), 2016 46th Annual IEEE/IFIP International Conference On*, pages 156–167. IEEE.
- [47] Leslie Lamport. The part-time parliament. 16(2):133–169.
- [48] Diego Ongaro and John Ousterhout. In search of an understandable consensus algorithm. In *2014 USENIX Annual Technical Conference (USENIX ATC 14)*, pages 305–319.
- [49] Fred B. Schneider. Byzantine Generals in Action: Implementing Fail-stop Processors. 2(2):145–154.
- [50] Leslie Lamport, Dahlia Malkhi, and Lidong Zhou. Vertical paxos and primary-backup replication. In *Proceedings of the 28th ACM Symposium on Principles of Distributed Computing*, pages 312–313. ACM.
- [51] John MacCormick, Nick Murphy, Marc Najork, Chandramohan A. Thekkath, and Lidong Zhou. Boxwood: Abstractions as the Foundation for Storage Infrastructure. In *OSDI*, volume 4, pages 8–8.
- [52] John MacCormick, Chandramohan A. Thekkath, Marcus Jager, Kristof Roomp, Lidong Zhou, and Ryan Peterson. Niobe: A practical replication protocol. 3(4):1.
- [53] David G Andersen, Jason Franklin, Michael Kaminsky, Amar Phanishayee, Lawrence Tan, and Vijay Vasudevan. FAWN: A fast array of wimpy nodes. In *Proceedings of the ACM SIGOPS 22nd Symposium on Operating Systems Principles*, pages 1–14. ACM.



- [54] Mahesh Balakrishnan, Dahlia Malkhi, Vijayan Prabhakaran, Ted Wobber, Michael Wei, and John D Davis. CORFU: A Shared Log Design for Flash Clusters. In *NSDI*, pages 1–14.
- [55] Michael Wei, Amy Tai, Christopher J Rossbach, Ittai Abraham, Maithem Munshed, Medhavi Dhawan, Jim Stabile, Udi Wieder, Scott Fritchie, Steven Swanson, and others. vCorfu: A Cloud-Scale Object Store on a Shared Log. In *NSDI*, pages 35–49.
- [56] Mahesh Balakrishnan, Dahlia Malkhi, Ted Wobber, Ming Wu, Vijayan Prabhakaran, Michael Wei, John D Davis, Sriram Rao, Tao Zou, and Aviad Zuck. Tango: Distributed data structures over a shared log. In *Proceedings of the Twenty-Fourth ACM Symposium on Operating Systems Principles*, pages 325–340. ACM.
- [57] Leslie Lamport. Paxos made simple. 32(4):18–25.
- [58] Leslie Lamport, Robert Shostak, and Marshall Pease. The Byzantine generals problem. 4(3):382–401.
- [59] Douglas B. Terry, Marvin M. Theimer, Karin Petersen, Alan J. Demers, Mike J. Spreitzer, and Carl H. Hauser. Managing Update Conflicts in Bayou, a Weakly Connected Replicated Storage System. 29.
- [60] Douglas B. Terry, Alan J. Demers, Karin Petersen, Mike J. Spreitzer, Marvin M. Theimer, and Brent B. Welch. Session guarantees for weakly consistent replicated data. In *Parallel and Distributed Information Systems, 1994., Proceedings of the Third International Conference On*, pages 140–149. IEEE.
- [61] Maurice P. Herlihy and Jeannette M. Wing. Linearizability: A correctness condition for concurrent objects. 12(3):463–492.
- [62] Hagit Attiya and Jennifer L. Welch. Sequential consistency versus linearizability. 12(2):91–122.
- [63] David K. Gifford. Information Storage in a Decentralized Computer System.
- [64] Leslie Lamport. Time, clocks, and the ordering of events in a distributed system. 21(7):558–565.
- [65] Barry M. Leiner, Vinton G. Cerf, David D. Clark, Robert E. Kahn, Leonard Kleinrock, Daniel C. Lynch, Jon Postel, Larry G. Roberts, and Stephen Wolff. A brief history of the Internet. 39(5):22–31.
- [66] Navin Budhiraja, Keith Marzullo, Fred B. Schneider, and Sam Toueg. The primary-backup approach. 2:199–216.
- [67] Sanjay Ghemawat, Howard Gobioff, and Shun-Tak Leung. The Google file system. In *ACM SIGOPS Operating Systems Review*, volume 37, pages 29–43. ACM.

- [68] Michael Stonebraker. The case for shared nothing. 9(1):4–9.
- [69] Google. Google RPC.
- [70] Tushar D. Chandra, Robert Griesemer, and Joshua Redstone. Paxos made live: An engineering perspective. In *Proceedings of the Twenty-Sixth Annual ACM Symposium on Principles of Distributed Computing*, pages 398–407. ACM.
- [71] Jon Gjengset <jon@thesquareplanet.com>. Students’ Guide to Raft :: Jon Gjengset.
- [72] Blake Mizerany, Qin Yicheng, and Li Xiang. Etd: Package raft.
- [73] Caio Oliveira, Lau Cheuk Lung, Hylson Netto, and Luciana Rech. Evaluating raft in docker on kubernetes. In *International Conference on Systems Science*, pages 123–130. Springer.
- [74] Heidi Howard, Malte Schwarzkopf, Anil Madhavapeddy, and Jon Crowcroft. Raft refloated: Do we have consensus? 49(1):12–21.
- [75] Lisa Glendenning, Ivan Beschastnikh, Arvind Krishnamurthy, and Thomas Anderson. Scalable consistency in Scatter. In *Proceedings of the Twenty-Third ACM Symposium on Operating Systems Principles*, pages 15–28. ACM.
- [76] Jim Gray, Pat Helland, Patrick O’Neil, and Dennis Shasha. The dangers of replication and a solution. In *ACM SIGMOD Record*, volume 25, pages 173–182. ACM.
- [77] H. Howard, D. Malkhi, and A. Spiegelman. Flexible Paxos: Quorum intersection revisited.
- [78] Leslie Lamport. Generalized consensus and Paxos.
- [79] Brian F. Cooper, Raghu Ramakrishnan, Utkarsh Srivastava, Adam Silberstein, Philip Bohannon, Hans-Arno Jacobsen, Nick Puz, Daniel Weaver, and Ramana Yerneni. PNUTS: Yahoo!’s hosted data serving platform. 1(2):1277–1288.
- [80] Wyatt Lloyd, Michael J. Freedman, Michael Kaminsky, and David G. Andersen. Don’t settle for eventual: Scalable causal consistency for wide-area storage with COPS. In *Proceedings of the Twenty-Third ACM Symposium on Operating Systems Principles*, pages 401–416. ACM.
- [81] Robbert Van Renesse and Fred B Schneider. Chain Replication for Supporting High Throughput and Availability. In *OSDI*, volume 4, pages 91–104.
- [82] Mike Burrows. The Chubby lock service for loosely-coupled distributed systems. In *Proceedings of the 7th Symposium on Operating Systems Design and Implementation*, pages 335–350. USENIX Association.

- [83] Patrick Hunt, Mahadev Konar, Flavio Paiva Junqueira, and Benjamin Reed. ZooKeeper: Wait-free Coordination for Internet-scale Systems. In *USENIX Annual Technical Conference*, volume 8, page 9.
- [84] Philip A Bernstein, Colin W Reid, and Sudipto Das. Hyder-A Transactional Record Manager for Shared Flash. In *CIDR*, volume 11, pages 9–12.